

**Understanding Collective Intelligence in Agent-Based
Systems: an Information-Theoretic Approach to the
Measurement and Comparison of Intelligence in Groups**

by

Nader Chmait



Thesis

Submitted by Nader Chmait

for fulfillment of the Requirements for the Degree of

Doctor of Philosophy (0190)

Main supervisor: Assoc. Prof. David L. Dowe

Associate supervisors: Prof. David G. Green and Dr. Yuan-Fang Li

**Faculty of Information Technology, Clayton
Monash University**

October, 2017

© Copyright

by

Nader Chmait

2017

Copyright notices:

Copyright © 2017 Nader Chmait

1. Under the Copyright Act 1968, this thesis must be used only under the normal conditions of scholarly fair dealing. In particular no results or conclusions should be extracted from it, nor should it be copied or closely paraphrased in whole or in part without the written consent of the author. Proper written acknowledgement should be made for any assistance obtained from this thesis.
2. I certify that I have made all reasonable efforts to secure copyright permissions for third-party content included in this thesis and have not knowingly added copyright content to my work without the owner's permission.

Contents

List of Tables	vii
List of Figures	viii
List of Abbreviations	x
Abstract	xi
Publications	xiv
Acknowledgments	xvi
1 Introduction	1
1.1 The Notion of Collective Intelligence	1
1.2 What Is Missing?	3
1.2.1 Measurement and comparison of intelligence	3
1.2.2 Predicting agent accuracy	3
1.2.3 Crossing boundaries between different cognitive systems	3
1.3 Goals of This Thesis	4
1.4 Main Contributions	5
1.5 My Approach	6
1.6 Thesis Outline	6
2 Collective Intelligence: A Literature Review	9
2.1 Overview	9
2.2 Collective Intelligence: Areas of Study	10
2.2.1 An artificial world sculpted by nature	10
2.2.2 Human societies and the rise of intelligence	14
2.3 Measuring Individual and Collective Intelligence	17
2.3.1 Evaluating non-human animals	17
2.3.2 Evaluating machines and artificial agents	19
2.3.3 Evaluating human intelligence	22
2.3.4 Individual versus collective intelligence	24

3	Factors of Collective Intelligence	29
3.1	Overview	29
3.2	Introduction	30
3.3	On Collective Intelligence	30
3.4	Evaluating Artificial Agents: A Short Background	32
3.5	Agent-Environment Framework	33
3.6	Intelligence Test	34
3.6.1	The Λ^* (Lambda star) environment	35
3.6.2	Algorithmic complexity	36
3.6.3	Search space complexity	39
3.6.4	Further thoughts on complexity	39
3.6.5	Intelligence score	40
3.7	Implementation Details and Experimental Protocol	40
3.7.1	Setup and test parameters	40
3.7.2	Modularity and code re-use	42
3.7.3	Implementing new agent behaviours	42
3.8	Agent Types and Behaviours	42
3.8.1	Local search agent	43
3.8.2	Reinforcement learning agents	43
3.8.3	Expert (oracle) agent	45
3.8.4	Random agent	45
3.9	Communication Protocols	46
3.9.1	Stigmergy or indirect communication	46
3.9.2	Implicit leadership through auctions and bidding	47
3.9.3	Imitating super-solver agents	47
3.9.4	Harnessing the wisdom of the crowd by information aggregation	48
3.10	Experimental Setup	48
3.11	Results and Discussion	49
3.11.1	Collectives outperform individuals	49
3.11.2	Communication and interaction protocol	50
3.11.3	Uncertainty in the environment	51
3.11.4	Number of agents in a group	54
3.11.5	Time and intelligence	55
3.11.6	Algorithmic complexity and intelligence	57
3.12	Organisational Behaviour	58
3.13	Alternative Environments and Further Considerations	61
3.13.1	Measuring multiagent coordination	61
3.13.2	Fitness landscapes	61
3.13.3	Further thoughts on robust intelligence tests	62
3.14	Conclusions and Future Work	63

4	Observation Communication and Intelligence	65
4.1	Overview	65
4.2	Introduction	66
4.3	Measuring Information	67
4.4	Agent Interaction Modes	69
4.5	Experimental Setup	70
4.6	Results and Discussion	70
4.6.1	Indirect communication	70
4.6.2	Direct communication	74
4.6.3	Imitation	75
4.7	Conclusion	78
5	Mathematical Predictive Modelling	83
5.1	Overview	83
5.2	Introduction and Background	84
5.3	Motivations and Major Contributions	85
5.4	Desirable Properties for Assessment	86
5.5	A Predictive Model of Agent Accuracy	87
5.6	Evaluating Inference Abilities Using the C-test	88
5.6.1	Structure of the test	89
5.6.2	Reward function	89
5.6.3	Inductive inference and choice of C-test	89
5.6.4	Measuring abilities	90
5.7	Predicting Agent Performances	92
5.7.1	Relationship between accuracy and task difficulty	93
5.7.2	Inferring task difficulty	95
5.8	Collective Accuracy of Cooperative Agents	95
5.9	Analysing Individual and Group Accuracies	96
5.9.1	Comparing agent collectives	98
5.9.2	A fictitious example	99
5.10	Conclusion	100
6	Network Science and Intelligence	103
6.1	Overview	103
6.2	What Lies at the Heart of Collective Intelligence?	104
6.3	Information Theory, Complexity and Intelligence	106
6.4	Intelligence as Coping With Complexity	109
6.4.1	Measuring and comparing intelligence across various cognitive systems	110
6.5	Network Theory and Intelligence	111
6.6	Integrating Network Science and Real World Disciplines	112
6.7	Comparing Graph Complexities	114
6.8	Sample Test Problems and Their Encoding	117
6.8.1	The scheduling problem	117

6.8.2	The travelling salesman problem	119
6.9	Agents as Networks	121
6.10	From Intelligence Tests to Other Spheres	123
6.11	Summary and Future Work	126
7	Conclusion and Future Work	129
7.1	Overall Contribution	130
7.2	Brief Summary and Main Outcomes	130
7.3	Limitations and Directions for Future Work	133
	Last Words	135

List of Tables

2.1	Characteristics of complex swarming behaviours.	11
2.2	Computer science meets biology: intelligent algorithms.	13
2.3	Factors and characteristics impacting the performance of human collectives.	16
3.1	Scores of artificial agents over the anYnt intelligence test.	50
6.1	Complexity measures of example sequence completion problems.	108
6.2	The exam timetabling problem: three examples.	118

List of Figures

2.1	Collective intelligence: areas of study.	10
2.2	Intelligence tests for evaluating non-human animals and insects.	18
2.3	Simulations of agents and collectives operating in artificial environments.	20
2.4	Intelligence test interface for evaluating humans.	23
2.5	Sample C-test sequences.	23
2.6	Raven’s Progressive Matrices.	23
2.7	Wildcat Wells: an online game for evaluating human groups.	23
3.1	Factors of collective intelligence.	31
3.2	The agent-environment framework.	34
3.3	An illustration of the Λ^* testing environment.	37
3.4	Encoding the movement pattern of special objects.	39
3.5	Agent behaviours: a simplified UML class diagram.	43
3.6	Initial group topologies.	48
3.7	A plot of the test scores appearing in Table 3.1.	51
3.8	Shift in effectiveness over different environment uncertainties.	53
3.9	Intelligence scores recorded across different numbers of agents.	54
3.10	Intelligence scores recorded across different test evaluation times.	56
3.11	Intelligence scores across different task complexities.	57
3.12	Group organisational structures.	59
3.13	Group effectiveness across different organisational structures.	60
3.14	Agents searching a fitness landscape: a simulation.	62
4.1	Observation and communication ranges.	68
4.2	Influence of indirect communication and observation on intelligence.	71
4.3	Indirect communication: variation in scores.	72
4.4	Indirect communication: whisker plot.	73
4.5	Influence of direct communication and observation on intelligence.	75
4.6	Direct communication: variation in scores.	76
4.7	Direct communication: whisker plot.	77
4.8	Influence of imitation and observation on intelligence.	78
4.9	Imitation: variation in scores.	79
4.10	Imitation: whisker plot.	80
4.11	Gradient difference of intelligence scores.	81

5.1	C-test intelligence scores.	93
5.2	Agent accuracy predicted by the IRT model.	94
5.3	Shift in accuracy across different problem settings.	94
5.4	Lower bound on accuracy guarantee.	95
5.5	Comparison of individual and collective accuracies.	97
5.6	Voting accuracy of three different team arrangements	100
6.1	Example testing environments.	106
6.2	Comparing intelligence between cognitive systems: Diagram I.	111
6.3	Comparing intelligence between cognitive systems: Diagram II.	112
6.4	Network structures and topologies.	112
6.5	Encoding graphs as binary strings.	113
6.6	Comparing intelligence between cognitive systems: Diagram III.	115
6.7	Different graph topologies and their complexity measures.	116
6.8	Topological graph complexity.	117
6.9	Modelling timetabling problems as graphs.	118
6.10	Comparing bumblebees to artificial agents.	119
6.11	Important group structures and templates.	125

List of Abbreviations

AI	Artificial Intelligence
AIT	Algorithmic Information Theory
anYnt	Anytime Universal Intelligence Test
IQ	Intelligence Quotient
IRT	Item Response Theory
\mathcal{LS}	Local Search
\mathcal{RL}	Reinforcement Learning
\mathcal{WOC}	Wisdom of the Crowd

Understanding Collective Intelligence in Agent-Based Systems: an Information-Theoretic Approach to the Measurement and Comparison of Intelligence in Groups

Nader Chmait

████████████████████
████████████████████
Monash University, 2017

Supervisor: Assoc. Prof. David L. Dowe

████████████████████
Associate Supervisors: Prof. David G. Green and Dr. Yuan-Fang Li
████████████████████

Abstract

Collective intelligence occurs in a wide range of areas such as social sciences, economics, biology and computer science. Famous applications include crowd-sourcing, public policy, recommendation systems, social computing, swarm intelligence and complex adaptive systems.

Despite the remarkable advancements in recent years, the vast majority of research on collective intelligence has investigated its emergence in isolation, that is, either within a limited range of disciplines, or at the level of one particular cognitive system. Thus, links are still missing to connect fundamental characteristics that are shared among these studies. Another serious limitation inhibiting our understanding of collective intelligence is the lack of *quantitative* analysis of intelligence with regards to task difficulty, and the comparison of individual agent performance to group performance. Therefore, the central question of this thesis is formulated as follows: What are the major characteristics and properties shaping the spread of intelligence *across various cognitive systems* and environments, and how to *quantitatively* measure and predict their influence on the performance of (individual and collectives of) agents? Consequently, this thesis aims to better understand the phenomenon of collective intelligence in various sorts of agent-based systems across the three cognitive systems: human, animal and machine.

There are several outcomes arising from this thesis. A new understanding of agent group performance and dynamics is established. This is achieved by showing how a range of factors and properties, that are inherent to agent groups of various cognitive types, quantitatively shape the agents' performance across different environment and problem settings. For instance, using formal intelligence tests, it is

discussed how (much) important factors, like *task information-theoretic complexity, the interaction mode between agents, their organisational structure, their observation and communication abilities and their decision-making dynamics*, bear influence on the agents' overall performance. Moreover, intrinsic dependencies between the examined factors are identified and measured.

While intelligence test scores are accurate measures of some abilities associated with the evaluated agent, they are an unreliable predictor of the agent's performance under different task difficulties or other problem settings. Therefore, a new mathematical predictive model is devised and used to predict the accuracy of agents over tasks of specific quantifiable complexities. The proposed model has several advantages. It makes it possible to avert the perpetual need to simulate agents over intelligence tests every time we need to predict their performance under a different problem configuration. A lower bound on agent accuracy can be guaranteed with respect to task complexity and the breadth of its solution space using the model. This in turn enables us to formulate the relationship between agent selection cost, task difficulty and accuracy as optimisation problems. Further results indicate the settings over which a group of agents can be more or less accurate than individual agents or other groups.

In the final parts of this thesis, I present a new perspective for comparing intelligence between non-uniform types of agents, operating in vastly different environments and contexts. Common grounds for evaluation are provided using a methodology for abstracting tasks and modelling environments as network graphs, showing how to measure their complexities. This is further used in the endeavour to connect studies of intelligence to other spheres, notably business decision-making and management.

Overall, this thesis provides general guidelines that give insight into how to explore the potential of collectives across different cognitive systems and research disciplines. It also provides initial forays towards bridging different research disciplines in which collective intelligence might occur, and consequently cross-fertilising diverse areas of study ranging from businesses and large organisations to social sciences and fundamental biology.

Understanding Collective Intelligence in Agent-Based Systems: an Information-Theoretic Approach to the Measurement and Comparison of Intelligence in Groups

Declaration

I declare that this thesis is my own work and has not been submitted in any form for another degree or diploma at any university or other institute of tertiary education. Information derived from the published and unpublished work of others has been acknowledged in the text and a list of references is given.

Nader Chmait
Saturday 14th October, 2017

Publications

List of publications arising from this thesis.

Published works

Nader Chmait, David L. Dowe, David, Yuan-Fang Li and G. Green (2017). *An Information Theoretic Predictive Model for the Accuracy of AI Agents Adopted From Psychometrics*, Proceedings of the 10th International Conference on Artificial General Intelligence (AGI), Vol. 10414 of Lecture Notes in Artificial Intelligence (LNAI), Chapter 21, Melbourne, Australia, Springer. [**Winner of the 2017 Kurzweil Best Paper Prize**].
https://link.springer.com/chapter/10.1007/978-3-319-63703-7_21

Jose Hernández-Orallo, Marco Baroni, Jordi Bieger, **Nader Chmait**, David L. Dowe, Katja Hofmann, Fernando Martínez-Plumed, Claes Strannegård, Kristinn R. Thórisson (2017). A New AI Evaluation Cosmos: Ready to Play the Game?, *AI Magazine, Association for the Advancement of Artificial Intelligence* **38**(3):66–69.
<https://www.aaai.org/ojs/index.php/aimagazine/article/view/2748>

Nader Chmait, David L. Dowe, Yuan-Fang Li, David G. Green, and Javier Insa-Cabrera (2016). *Factors of collective intelligence: How smart are agent collectives?*, Proceedings of the 22nd European Conference on Artificial Intelligence (ECAI), The Hague, The Netherlands, Vol. 285 of Frontiers in Artificial Intelligence and Applications, pp. 542–550.
<http://ebooks.iospress.nl/volumearticle/44798>

Nader Chmait, Yuan-Fang Li, David L. Dowe, and David G. Green (2016). *A dynamic intelligence test framework for evaluating AI agents*, Proceedings of the 1st International Workshop on Evaluating General-Purpose AI (EGPAI 2016), European Conference on Artificial Intelligence (ECAI) 2016, pp. 1–8.
<http://www.ecai2016.org/content/uploads/2016/08/W14-EGPAI-2016.pdf>

Nader Chmait, David L. Dowe, David G. Green, and Yuan-Fang Li (2015). *Observation, communication and intelligence in agent-based systems*, in J. Bieger, B. Goertzel and A. Potapov (eds), Proceedings of the 8th International Conference on Artificial General Intelligence, Berlin, Germany, Vol. 9205 of Lecture Notes in Artificial Intelligence (LNAI), Springer, pp. 50–59.
http://dx.doi.org/10.1007/978-3-319-21365-1_6

Other manuscripts

Nader Chmait, David L. Dowe, David G. Green and Yuan-Fang Li (2017). *Coping with complexity: a multi-disciplinary survey on collective intelligence and its measurement*, Currently under review.

Nader Chmait, David L. Dowe, David G. Green, Yuan-Fang Li, and Javier Insa-Cabrera (2015). *Measuring universal intelligence in agent-based systems using the Anytime Intelligence Test*, Technical Report 2015/279, FIT, Clayton, Monash University, Australia, p. 13.
<http://www.csse.monash.edu.au/publications/2015/tr-2015-279-full.pdf>

Permanent Address: Faculty of Information Technology, Clayton
Monash University
Australia

This thesis was typeset with $\text{\LaTeX} 2_{\epsilon}$ ¹ by the author.

¹ $\text{\LaTeX} 2_{\epsilon}$ is an extension of \LaTeX . \LaTeX is a collection of macros for \TeX . \TeX is a trademark of the American Mathematical Society. The macros used in formatting this thesis were written by Glenn Maughan and modified by Dean Thompson, David Squire and the thesis author Nader Chmait of Monash University.

Acknowledgments

This PhD journey has been the most challenging and at the same time the most exciting experience in my life so far. At the end of this journey, I would like to express my appreciation to my academic supervisors Assoc. Prof. David L. Dowe, Prof. David G. Green and Dr. Yuan-Fang Li for their guidance throughout my candidature. Both Davids and Yuan-Fang were a great source of inspiration for me.

I would like to thank Prof. José Hernández-Orallo who invited me for a study-away at the Universitat Politècnica de València. José also offered me the opportunity to join him and others in organising the second *International Workshop on Evaluating General-Purpose AI* (Chmait, Hernández-Orallo, Martínez-Plumed, Strannegård and Thórisson, 2017) which was a joyful and fulfilling experience. Working with José has also helped me to connect with other researchers and professionals in my area of research. My gratitude also goes to Dr. Javier Insa-Cabrera who welcomed me and accompanied me during my stay in València.

I wish to acknowledge the financial support I received from the Faculty of Information Technology as part of the Postgraduate Publication Award scheme. I also wish to acknowledge the efforts of the academics and staff members of the Faculty of Information Technology at Monash University, especially my milestones' review panel members Prof. Alan Dorin, Prof. Bala Srinivasan, and Dr. Julian Garcia for their valuable feedback, Dr. Jan Carlo Barca who welcomed me to the Monash Robotics Swarm Lab and with whom I discussed potential applications of my research in the fields of robotics and swarm intelligence, as well as the former and current graduate student services' coordinators Helen Cridland, Aidan Solla and Danette Deriane who supported and encouraged me across my PhD research milestones.

Finally, I would like to thank my lovely family, colleagues and friends for their tremendous support during my whole candidature. A special thanks goes to Tenindra Abeywickrama, Ehsan Shareghi, Srinibas Swain and Thomas Hendrey who enthusiastically listened to me every time we discussed my research together.

Without you all, this thesis would not have been possible. It has been such an exceptional journey of exploration, hard work and self discovery.

Nader Chmait

Monash University
October 2017

Chapter 1

Introduction

Aucun de nous ne sait ce que nous savons tous, ensemble.

—Euripide, Dramaturge

1.1 The Notion of Collective Intelligence

Studying the major characteristics and factors shaping the spread of *intelligence* in human and non-human *collectives* has several advantages and a wide range of implications. For example, it is central for the understanding of the social behaviours of these collectives, as well as, among other things, enhancing business decision-making and improving management strategies by harnessing the potential of the members of the collective.

Collective intelligence (Weschler, 1971) is an old concept that has been studied at least as far back as the mid-1780s as an interpretation of Condorcet’s jury theorem (De Caritat Marquis De Condorcet, 1785). Later adopted by many scientists and philosophers, the term “collective intelligence” has come to express intelligent behaviour in groups. A large number of definitions of collective intelligence now exist and vary in nature and scope. For instance, collective intelligence has been described as a property emerging from the synergies among different entities (Glenn et al., 2014, p. 7), “a form of universally distributed intelligence, constantly enhanced, coordinated in real time, and resulting in the effective mobilisation of skills” (Lévy, 1997, p. 13), or further as the “capability for a group of people to collaborate in order to decide upon its own future and reach it in a complex context” (Noubel, 2004, p. 19), and lately – using a more general definition – as any “group of individuals doing things collectively that seem intelligent” (Malone, 2008).

In broad terms, one can simply interpret the term collective intelligence as any form of *intelligence* that is exhibited by a *collective*. A collective is a (distributed) system of agents where each agent² (referring to an autonomous entity, subject or component of that system) has a private utility function it is trying to maximise, along with a world utility function that measures the overall system’s performance (Tumer and Wolpert, 2004). In

²The notion of an agent and its different types and classifications (biological, robotic or computational) has been studied in the context of intelligence in (Pfeifer and Scheier, 2001, Chap. 1).

this thesis the term *collective* is used in the broader context as an abstraction of any, large or small, group of interactive agents sharing some common objective.

The term *intelligence* itself has been given many definitions. As quoted in (Gregory and Zangwill, 1987, p. 376), R. J. Sternberg comments that “viewed narrowly, there seem to be almost as many definitions of intelligence as there were experts asked to define it”. For example, intelligence was defined by Weschsler as the “composed or global ability of an individual to act purposeful, think reasonably, and to effectively deal with its environment” (Wechsler and Hardesty, 1964, p. 13) and more recently as the “ability to learn, to understand, and to adapt to an environment by using own knowledge” (Leimeister, 2010, Section 2). Legg and Hutter overviewed in (Legg and Hutter, 2007, Sec. 2.6 and 2.7) more than ten different definitions of intelligence and identified commonalities between them, showing that intelligence is recurrently “seen as a property of an individual who is interacting with an external environment” and it is related to the individual’s ability to succeed or *profit*.

There is evidence of the phenomenon of collective intelligence in all sorts of groups, whether human, animal or machine (Malone and Bernstein, 2015; Leimeister, 2010). Understanding (the dynamics behind) the intelligent behaviour among distributed interactive agents and groups can be highly beneficial. New models and promising solutions to a variety of (e.g., multiagent, optimisation and prediction-like) problems are now possible by harnessing the *wisdom of the crowd* (Surowiecki, 2005). The motivations are numerous. For instance, collective intelligence can enhance business operations and transform modern enterprise systems by allowing for personalised user content and new forms of user interactions (such as collaborative filtering) (Nagalakshmi and Joglekar, 2011). These characteristics have increased the competitive advantage of many (online) businesses such as Facebook, Amazon and eBay. Apart from escalating business profits, collective intelligence is now implicit in the design of robotics and artificial agents (Beni, 2004; Bonabeau et al., 1999). New models for adequately solving hard computational problems are devised by understanding how simple individual actions can lead to complex social organisations, and to intelligent behaviour on a larger scale. Different types of collectives have been studied and evidence for intelligent collective behaviour has been found in each case. Examples from various fields and disciplines are given below:

- (a) biology (Garnier et al., 2007; Conradt and Roper, 2005; Millonas, 1994), e.g., insect colonies, fish schooling, bird flocking, microbial behaviour,
- (b) social sciences, psychology and economy (Benkler, 2006; Weschsler, 1971; Camerer et al., 2011) e.g., human societies and polities, group (social) behaviour, finance, efficient markets, crowd-sourcing, human-computation, and
- (c) computer science (Bonabeau et al., 1999; Engelbrecht, 2006; Yang, 2010a; Wolpert et al., 2013) e.g., distributed artificial intelligence, swarm intelligence, artificial life, nature-inspired and evolutionary computation.

Recently, the *Center of Collective Intelligence* at MIT also reported on a large number of interesting studies (Malone and Bernstein, 2015) that were conducted in many of the above

areas. This diversity of studies highlights the importance of investigating the concept of collective intelligence as it is repeatedly witnessed throughout different types of interactive entities, systems, and disciplines.

1.2 What Is Missing?

Interest in the study of collective intelligence is growing rapidly³. Numerous works have investigated this topic in one or more disciplines. Yet, major gaps still exist in this area of study. Despite the extant literature on collective intelligence, important questions like “how does the effectiveness of a group **quantitatively** compare to that of its isolated members?” and “are there some general rules or properties shaping the spread of **intelligence across various cognitive systems** and environments?” remain somewhat of a mystery. Major challenges stem from the previous questions. I describe these challenges in the following sections.

1.2.1 Measurement and comparison of intelligence

The first challenge is measuring the effectiveness of both individual agents and groups of interactive agents, and comparing them, across different problem/environmental settings and complexities. The second challenge is to identify and analyse the underlying factors that may have shaped (e.g., lead to any observed differences in) the effectiveness of the evaluated agents. Questions like “how (much) does the number of agents in the group influence its performance?” illustrate part of these challenges.

1.2.2 Predicting agent accuracy

Intelligence tests are perhaps the most efficient and frequently used tools for measuring the average performance of a (human, animal or artificial) subject over a set of tasks or environments. Such tests are substantial whenever two subjects or systems need to be assessed and their performances compared, and hence make powerful tools to tackle the questions I raised in the first paragraph of Section 1.2. Nevertheless, intelligence tests might fall short when it comes to predicting the accuracy of a subject over a particular task (or environment) *complexity* without actually administering that task to the testee, and thus a different approach must be sought for this purpose. In other words, intelligence tests are good indicators of the overall performance but are not (rigorous or specific enough to serve as) accurate predictive models.

1.2.3 Crossing boundaries between different cognitive systems

Another major challenge regarding the assessment of intelligence is how to practically compare intelligence in different contexts, e.g., among agents that cannot operate in common

³Google Books Ngram Viewer (Michel et al., 2011) shows that the use of the term “collective intelligence” in books has risen significantly based on the *2009 English One Million corpus* data. Search queries on Google Scholar with the keyword “collective intelligence” in the title return only 16 results for all articles published between 1950 and 1990, and 1810 results for those published since 1990.

environments. This gap has led to the widespread disparity and inconsistency in defining and measuring intelligence (Hernández-Orallo et al., 2016; Fulker and Eysenck, 2012; Hernández-Orallo and Dowe, 2010)(Dowe, 2013, Sec. 4.4), inter and even intra-disciplines, which obviously extends to the measurement of collective intelligence. Consequently, it is very difficult to cross the boundaries between human, animal and machine entities; each time we need to analyse the manifestation of intelligence amongst their groups of individuals. This significantly inhibits the extent to which we can exploit the outcomes from one field in order to cross-fertilise others.

1.3 Goals of This Thesis

Given the above considerations, this thesis project has three main goals ((**G01**), (**G02**) and (**G03**)) that are outlined below. Each goal is broken down into smaller objectives (denoted by **Obj** followed by a numerical identifier).

(G01): The first main goal is to measure and quantitatively compare the agents' individual and collective performances. This goal consists of three objectives.

Obj01: In order to measure performance I will implement and use a formal (general and dynamic) intelligence test to *quantify* the effectiveness of the evaluated agents.

Obj02: In order to compare performance between groups and individuals I will identify some of the main factors influencing the behaviour of (different types of) groups and individual agents. These factors can either be associated with the testing environment or the structure and characteristics of the evaluated agents.

Obj03: I will analyse how (much) the identified factors affect the performance of the agents, and how these factors are related to one another. This objective entails measuring the (individual and simultaneous) impact of these factors on performance, and revealing the dependencies between them.

(G02): The second main goal is to design a mathematical predictive model to help predict the accuracy of agents across different problem settings and difficulties.

Obj04: Measure the ability of an agent over a certain class of tasks. This involves interpreting experimental outcomes in a similar manner as in *Obj01*.

Obj05: Develop a model that predicts the accuracy of an agent of some measured ability over different problem difficulties/settings without the need to evaluate the agents over these problems.

Obj06: Use the model to compare and analyse the (predictive) accuracies of individual agents and groups.

(G03): The third main goal is to cross boundaries between cognitive systems and provide a methodology to compare intelligence among agents that are (biologically

or physically) unfit to operate in common environments. While this might be an ambitious goal, I aim to propose a preliminary infrastructure that provides some initial steps towards comparing intelligence between agents (belonging to different cognitive systems and) operating in different contexts.

Obj07: Use network theory to model/abstract different types of environments.

Obj08: Present a technique to measure the *complexities* of such environments.

Obj09: Provide a series of steps that can be used to compare intelligence between agents operating (under different configurations) in such environments.

1.4 Main Contributions

There are three main contributions in this thesis that fill in the gaps in the literature identified earlier in Section 1.2. These contributions are described below.

Contribution One: A new understanding of agent group performance and dynamics is established. This is achieved by showing how a range of factors and properties, that are inherent to agent groups of various cognitive types, quantitatively shape the agents' performance across different environment and problem settings.

Contribution two: Intelligence tests are accurate measures of average agent performance but are an unreliable predictor of the agent's performance under different task difficulties or other problem settings. This limitation is overcome by devising a novel mathematical predictive model for approximating agent (group) accuracy over problems of different complexities without the perpetual need to administer these problems to the agents.

Contribution three: Quantitative comparison of performance between diverse kinds of agents operating in substantially different environments and contexts is made possible by providing common grounds for evaluation. Studies of intelligence are also connected to other spheres, notably business decision-making and management. This is accomplished by presenting a methodology to abstract tasks and model environments as network graphs and measure their complexities prior to evaluation.

Looking at the big picture, these contributions have many implications for the design of intelligent multiagent systems, the understanding of their *social* behaviour, and the prediction of their capacity for intelligence. These contributions serve in turn as general guidelines that give insight into how to explore the potential of collectives across different cognitive systems and research disciplines. A detailed description of the research outcomes that have led to the above-described contributions is given in subsequent chapters. It is further discussed how these outcomes are relevant to many important real-world problems at the level of organisations and teams, and how they might enhance our understanding of different nature-inspired and swarm-like behaviours.

1.5 My Approach

In order to explore the characteristics of collective intelligence across different research areas and cognitive systems we need to look at this phenomenon and analyse it from different, high- to low-level, perspectives.

To achieve the above, I begin by surveying the latest research on collective intelligence and its measurement conducted across diverse cognitive systems and disciplines and identify common threads that show how they relate to each other.

I continue by empirically assessing artificial agents over formal information-theoretic intelligence tests from the literature. I describe how to quantify the complexity of the testing environment and its assessment tasks. I simulate the behaviours of such agents, both in isolation and collectively, in order to investigate whether groups can lead to more intelligent systems. I evaluate diverse ways in which these agents can be put together in one group by exploring, among other things, different collective decision-making techniques, organisational structures and communication protocols. This experimental approach will be used to compare the agents' performances over the different evaluation settings and to quantify the influence of the examined factors that demonstrated an impact on their scores.

In addition to such experiments, I perform some mathematical predictive modelling in order to overcome some of the limitations of simulations and intelligence tests. For instance, I propose a formal model to predict the accuracy of a (group of) agents(s) across problems (or tasks) of well-defined complexities without the perpetual need to administer the agents to these problems.

In the final part of this thesis I address the challenge of (measuring and) comparing the intelligence of agents operating in substantially different environments. Due to this challenge, it seems absurd to try to make sense of any results collected from testing over such environments given the large differences and non-uniformity between them. Thus, I start by identifying some general properties relevant to the quantification of intelligence that transcend any particular cognitive system and discipline. I show how we can make quantitative comparisons between the effectiveness of a collective, its isolated members, and other kinds of collectives. Furthermore, I describe some of the obstacles that hinder the assessment of collective intelligence which may be either inflicted by the examined cognitive system type, or the environment itself. Finally, I present a preliminary step towards solving this problem. I propose a methodology to model vastly different environments and problems as network graphs and measure their complexities. I also give examples of how to quantitatively compare the effectiveness of agents across these environments.

1.6 Thesis Outline

This thesis is divided into seven chapters. The outcomes from these chapters incrementally contribute to the fulfilment of my identified goals as described below.

Chapter 2

This chapter gives a detailed historical background on the notion of collective intelligence and its measurement, introducing many of its key application and research areas. Chapter 2 also raises questions that are fundamental to the understanding of collective intelligence and the identification of its underlying characteristics, many of which I revisit in subsequent chapters. The scope goes beyond investigating a single cognitive system and discusses the literature of the measurement of intelligence in humans, animals and machines. This chapter also discusses some of the difficulties and unsolved problems of comparing intelligence between different types of systems and makes an overture to the forthcoming chapters.

Chapter 3

This chapter addresses (G01). A number of factors influencing and hindering the collective intelligence of interactive cognitive systems are identified and studied in this chapter. After conducting a series of controlled experiments over formal information-theoretic intelligence tests from the literature, I measure the (independent and simultaneous) influence of the examined factors on intelligence. Furthermore, I investigate how the organisational, or network, structure of equally sized groups shapes their effectiveness.

Chapter 4

Chapter 4 also addresses (G01). It mainly tackles objectives *Obj02* and *Obj03*. Here I focus on two main factors known to influence the performance of multiagent systems and their capacity for intelligence. These factors are the communication and observation abilities of their agents. I empirically measure and compare the effectiveness of cooperative agents of different observation/perception and communication abilities and highlight the circumstances under which they achieve optimal performance. I also discuss the dependency between the studied factors across different agent cooperation scenarios and interaction modes.

Chapter 5

Chapter 5 addresses (G02). In the previous chapters I principally adopted an experimental approach to achieve our purposes, whereas in this chapter I introduce a new mathematical model to quantitatively estimate the accuracy of artificial agents. I derive my proposed model by introducing notions from algorithmic information theory into Item Response Theory (IRT), which is a well-known (psychometric) measurement paradigm. I demonstrate the model by predicting the accuracy of isolated and cooperative artificial agents over inductive inference problems of varying complexities and (breadth of) solution spaces. Furthermore, I indicate the settings over which a group of agents can be more or less accurate than individual agents or other groups when solving cognitive tasks.

Chapter 6

Chapter 6 addresses (G03). In this chapter I link the theory of complexity to some real world problems, and show how the results from the science of complexity can be used to cross-fertilise the different research areas relevant to collective intelligence. I present a methodology that might be used to model a large number of (natural and artificial) environments and tasks as networks of different structures. I then show how, using the presented methodology, one can cross boundaries between human, animal and artificial cognitive systems and how to compare their performances without the need for a universal intelligence test.

Chapter 7

In this chapter I combine many of the conclusions drawn in earlier chapters and summarise the main contributions of this thesis. A short overview of some of the main outcomes resulting from chapters 2 to 6 is given below:

- *Chapter 2:* A cross-disciplinary survey of a wide range of studies conducted on collective intelligence and its measurement.
- *Chapter 3:* Quantitative analysis of agent individual and group performances over different environmental and cooperative settings.
- *Chapter 4:* Quantitative analysis of the trade-off between agents' communication and observation capacities on overall performance and their dependencies.
- *Chapter 5:* Development of a mathematical model for predicting the accuracy of agents under different problem complexities.
- *Chapter 6:* An approach for the measurement and comparison of intelligence between different cognitive systems operating in substantially different environments.

Chapter 7 also discusses some of the limitations of this study, gives some directions for future work, and raises many state-of-the-art and open questions in artificial intelligence.

Chapter 2

Collective Intelligence: A Literature Review

Order is created from chaos; patterns are revealed; and systems are free to work out their errors and problems at their own level. What natural systems can teach humanity is truly amazing.

—L. K. Samuels, In Defense of Chaos (2013)

The research in this chapter will provisionally be published in the following articles:

- Nader Chmait, David L. Dowe, David G. Green and Yuan-Fang Li. (2017). *Coping with complexity: a multi-disciplinary survey on collective intelligence and its measurement*, To be submitted to the AI Journal, Elsevier.
- Jose Hernández-Orallo, Marco Baroni, Jordi Bieger, Nader Chmait, David L. Dowe, Katja Hofmann, Fernando Martínez-Plumed, Claes Strannegård, Kristinn R. Thórisson (2017). A New AI Evaluation Cosmos: Ready to Play the Game?, *AI Magazine, Association for the Advancement of Artificial Intelligence* **38**(3):66–69.

2.1 Overview

Collective intelligence is observed in both the natural and artificial worlds. It has been studied repeatedly in many fields and disciplines, from life sciences and biology (animal herds and insect colonies), to social sciences and psychology (human societies, polities, and organisations harnessing the wisdom of the crowd), and computer science (artificial life and nature-inspired evolutionary computation).

In this chapter I survey the latest research on collective intelligence and its measurement that was conducted across diverse cognitive systems and disciplines. Some of the surveyed works will be revisited in later chapters where I identify common threads of how they relate to one another, and show how the intelligence of a collective compares to that of its members and other types of collectives.

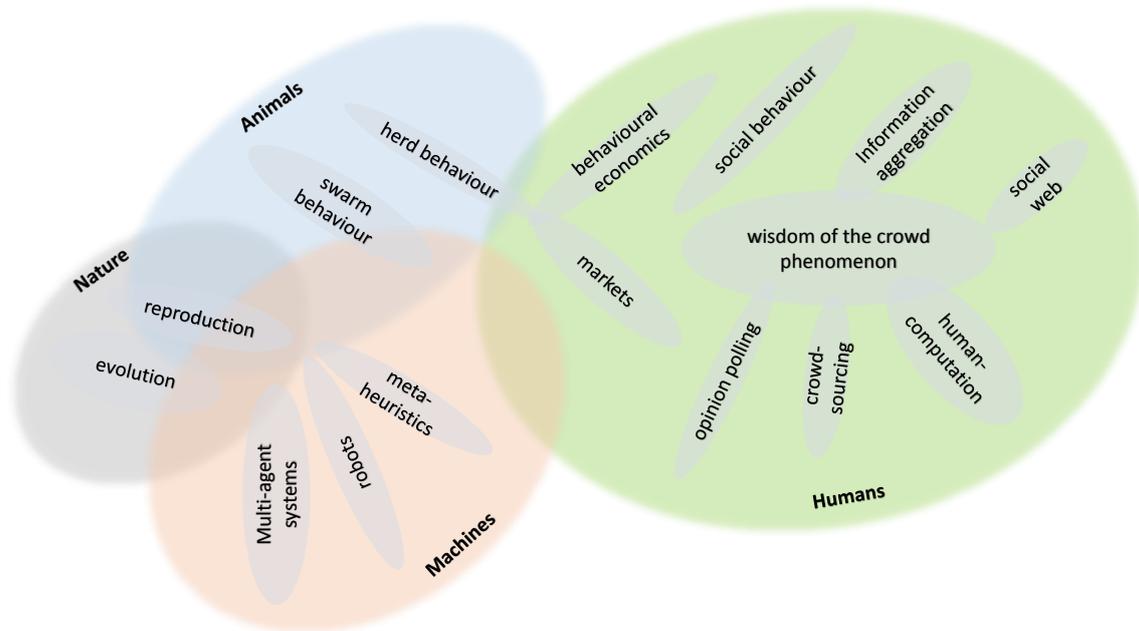


Figure 2.1: Areas of study of collective intelligence (and some of their intersections) that I aim to survey in this chapter.

2.2 Collective Intelligence: Areas of Study

I will first look (in the next section) at the different areas of study relevant to the notion of collective intelligence, many of which are depicted in Figure 2.1. I will attempt to separate these studies into two main subsections: subsection 2.2.1, where I investigate collective intelligence both in nature and in the animal kingdom, and its large influence on the Artificial Intelligence (AI) community, and subsection 2.2.2, in which I focus on human collectives, and the emergence of intelligence amongst human groups and societies throughout a large number of applications and disciplines. I continue (in Section 2.3) by reviewing some of the main techniques for evaluating and measuring individual and collective intelligence, in all three cognitive systems: human, animal and machine, as well as hybrid systems.

2.2.1 An artificial world sculpted by nature

Animals and insects are remarkable social entities. Driven by the need to exploit their environment for food and shelter, or to escape predators and avoid injury, they have developed exceptional cognitive and social abilities (Couzin, 2009). Their emergent behaviour (Bonabeau et al., 1999; Garnier et al., 2007), often referred to as swarm intelligence, heavily inspired artificial intelligence (Parpinelli and Lopes, 2011). Different research studies (Garnier et al., 2007; Conradt and Roper, 2005; Millonas, 1994; Bonabeau et al., 1999; Miller et al., 2013; Reid et al., 2015) have suggested principles that may be responsible for the intelligent and complex behaviour in non-human animal groups, and swarms in particular. Most of these principles reduce down to the ability of self-organisation, randomness and adaptability, as well as the interaction between the members

of the collective. I list in Table 2.1 a few of those mechanisms underlying self-organisation and the complex swarming behaviours following the works in (Bonabeau et al., 1999; Millonas, 1994; Garnier et al., 2007; Reynolds, 1987).

<i>Characteristic name</i>	<i>Description</i>
1) Communication type	<p><i>Stigmergy</i> (Grassé, 1959), a form of indirect communication which consists of modifying the environment (serving as a shared state memory) and exploiting the asynchronous response from different individuals echoing the changes in their neighbourhood</p> <p><i>Direct communication</i> between the individuals through trophallaxis (direct contact), antennation or even reproduction. Another form of direct communication is broadcasting by propagating signals throughout the environment.</p> <p><i>Feedback</i>—positive or negative—from the individuals in the collective. These are different rules and interaction mechanisms promoting the creation of structures and stability inside the collective.</p>
2) Randomness	Fluctuations from the received feedback. Randomness leads to a somewhat stochastic decision-making or behaviour which happens to be very useful for enhancing search strategies and foraging by reducing the odds of getting trapped in local optima.
3) Emergence	Exhibition of emergent properties through local (and multiple) interactions among the members of the collective. Properties like cohesion, separation, and alignment relying on different (nearest neighbour, position centring, velocity matching, etc.) rules shape the behaviour of various animal collectives.
4) Stability	Resulting from the initial condition in the environment, the organisation of the collective and the random fluctuations.
5) Adaptability	The ability to efficiently switch behaviour (or update the fitness of the individuals) under the right and suitable conditions by evaluating the reward sent back by the environment or other external feedback.

Table 2.1: Characteristics of complex swarming behaviours: some mechanisms underlying self-organisation and the complex swarming behaviours following the works in (Bonabeau et al., 1999; Millonas, 1994; Garnier et al., 2007; Reynolds, 1987).

The above principles arose from the study of various kinds of entities in the animal kingdom and nature in general such as ant colonies, flocks of birds, schools of fish, foraging bees, bacteria, self-assembling proteins, glow-worms, fireflies, slime moulds, cockroaches, mosquitoes and many other organisms. Nonetheless, these principles are not restricted to animal swarms, but rather consist of general properties which can emerge from (the interactions within) any large modular reconfigurable system. For instance, self-organisation and other properties of emergence were studied in the multiagent systems community (Serugendo et al., 2006; Aliu et al., 2013; Bernon et al., 2006). The different types of communication protocols used in swarms such as broadcasting and information propagation throughout the environment (e.g., bees' waggle dance, fireflies' glow

intensity), direct communication (e.g., antennation, reproduction, trophallaxis by food or liquid exchange), and stigmergy (ant pheromone trails, nest building in social wasps or termites), have provided them with exceptional abilities to cooperatively solve complex problems. Not surprisingly, the multiagent system community has adopted many of these techniques and communication protocols in their approach to the design of self organisation, as for example by using direct interaction and communication (Mamei et al., 2004), as well as stigmergy (Karuna et al., 2004; Bourjot et al., 2003). Other approaches consisted of using artificial systems with dynamically adaptive agents (Gleizes et al., 1999; Maes, 1993), or even reinforcement (learning) mechanisms with predefined system architectures (Maturana and Norrie, 1996). For instance, a collective intelligence framework (COIN) (Wolpert and Tumer, 1999; Wolpert, 2004; Wolpert et al., 2013) was designed to engineer an agent's private reward utility function in such a way as to enhance the emergence properties of the collective. This was done by formally exploring the conditions sufficient for the emergent behaviour for collectives of independent reinforcement learning (Watkins and Dayan, 1992) agents.

To give some insight into how different species in the animal kingdom collectively solve complex problems, let's look again into how animal behaviours inspired the computer science community as most of these behaviours were used in the design and development of different algorithms and heuristics for solving hard, real-world (optimisation) problems. Table 2.2.1 summarises some of the most famous meta-heuristics, swarm and nature-inspired algorithms and gives a brief description of their origins.

The list in Table 2.2.1 is non-exhaustive, and the applications are numerous (Bonabeau et al., 1999; Bonabeau and Meyer, 2001; Kennedy et al., 2001; Eberhart and Shi, 2001; Engelbrecht, 2006; Ducatelle et al., 2010). Many other examples are available illustrating how evolution (Draves, 2008; Angeline, 1995), co-evolution (Van Veldhuizen and Lamont, 2000; Potter and De Jong, 2000), and more generally nature (Popkin, 2016; Yang, 2010a), gave birth to diverse meta-heuristics for solving many complex problems, and enhanced the design of cooperative multiagent learning (Panait and Luke, 2005). Note here the difference in approach to optimisation between some traditional AI algorithms, like the A* algorithm (Hart et al., 1968) on one hand, and nature-inspired algorithms on the other hand. While an A* algorithm using an admissible heuristic is optimal (e.g., finds shortest path), most nature-inspired algorithms don't guarantee optimality, but instead provide a sub-optimal, or adequate, solution to complex problems in a reasonable period of time.

The exploitation of animal and nature-inspired intelligence extends beyond software agents and heuristics to reach the physical world. For instance, biologically inspired robots (Fong et al., 2003; Beni, 2004) internally mimicking the social behaviour in living creatures were designed and manufactured, in many cases allowing for the opportunity to interact and cooperate with human beings by engaging in their daily activities. Similarly to artificial agents, the design of collective robots was mainly inspired by nature. Interactive robots however have one major distinction from artificial agents as they are embodied in their environment, in the sense that, they can perturb or get perturbed by it (Brooks, 1991). This distinction is very important because it might redefine the nature

<i>Algorithm or meta-heuristic</i>	<i>Origin and inspiration</i>
Ant Colony Optimization (Bonabeau et al., 1999; Dorigo et al., 2006; Dorigo and Stützle, 2009)	Ants foraging for food
Artificial Bee Colony (Karaboga, 2005)	Foraging in employed and onlooker bee phases
Bee Algorithm (Pham et al., 2006)	Bees foraging for food
Bee Hive (Wedde et al., 2004)	Bee organisation and communication
Marriage in Honey-bees Optimisation Algorithm (Abbass, 2001)	Bee mating/search for queen
Particle Swarm Optimization (Kennedy, 2011)	Swarm social behaviour, movement and organisation
Bacterial Foraging Algorithm (Passino, 2002)	Foraging/reproducing bacteria
Glow-worm Swarm Optimisation Algorithm (Passino, 2002)	Firefly bioluminescence/flashing while sensing neighbourhood environment
Slime Mould Optimisation Algorithm (Monismith and Mayfield, 2008)	Foraging and dispersal of Amoebae cells or organisms
Roach Infestation Optimisation Algorithm (Havens et al., 2008)	Cockroaches foraging for food
Cuckoo Search (Yang and Deb, 2009)	Brooding behaviour of cuckoos
Bat Algorithm (Yang, 2010b)	Echolocation strategies in bats
Firefly Algorithm (Yang, 2009)	Flashing fireflies movement and attractiveness
Social Spider Optimization algorithm (Cuevas et al., 2013)	Operational principles from the social-spider colony
Gravitational Search Algorithm (Rashedi et al., 2009)	Gravity and mass interaction physical laws
River Formation Dynamics (Rabanal et al., 2007)	River formation by eroding water
Intelligent Water Drops (Shah-Hosseini, 2009)	River paths flow
Altruism (Foster et al., 2006)	Altruistic behaviour according to Hamilton's rule of kin (Hamilton, 1964a,b) selection
Artificial Immune System (De Castro and Timmis, 2002)	Functionality of the immune system

Table 2.2: Computer science meets biology: a brief description of meta-heuristics, swarm and nature-inspired algorithms (pertinent to the notion of collective intelligence) that have been exploited in both disciplines of computer science and biology.

of (distributed) computation (and consequently its evaluation) to include the influence of environment (Hoffmann and Pfeifer, 2012), and also capture the integrated signals received from different *sensory systems* and the brain.

Despite all their differences, it is intriguing how the collective behaviours of so many different species provide us with effective strategies to efficiently solve a wide range of unpredictable and complex problems. As we will see in the next section some animal group behaviours, such as herd behaviour (Banerjee, 1992), are also shared with human animals and manifest at different levels of their individual decision-making process leading to complex social behaviours.

2.2.2 Human societies and the rise of intelligence

One of the most famous examples of collective intelligence as a tool for efficient prediction and estimation was Galton's 1906 experiment (Galton, 1907) of *guessing the weight of an Ox*, which emphasises the "wisdom of the crowd" phenomenon. Several real-life examples of the latter phenomenon were presented and analysed by Surowiecki (Surowiecki, 2005) in 2005, such as the *jelly beans in the jar* experiment, and the *Who wants to be a millionaire?* television show, in which the group's estimate is always superior to the vast majority of the individual guesses. In fact, using the crowd as a forecast tool by polling the collective opinion is widely used nowadays as part of a large number of (collaborative) crowd-sourcing, crowd-mining and other human-computation applications (Howe, 2008; Tapscott and Williams, 2008), which in turn have been exploited in diverse areas like economics, law, psychology, sociology, as well as politics. For instance, there are some famous resources and applications like *Amazon Mechanical Turk*, *Wikipedia*, *TripAdvisor*, and others (Burbank et al., 2011; Raddick et al., 2010), all of which harness the collective (and often specialised) knowledge of the crowds, and are used as training data for information retrieval (Denoyer and Gallinari, 2006), mining consumer behaviours (O'Connor, 2008) and other prediction purposes (Berinsky et al., 2012; Buhrmester et al., 2011; Malone et al., 2010). These applications have had a huge success and economical impact, in spite of their non-traditional design and modular characteristics, in which different units and (specialised) opinions are grouped together and recombined in order to design larger, richer and more accurate systems.

As the backbone of economy, the importance and impact of markets on human societies is immense. In the last two decades, numerous studies have been undertaken in behavioural economics (Camerer and Fehr, 2006; Schelling, 2006; Camerer et al., 2011), psychology and prediction markets. The studies investigated the consequences and advantages of information polling and information aggregation in estimating probabilities and making predictions (Wolfers and Zitzewitz, 2004; Benkler, 2006; Easley and Kleinberg, 2010), and they have shown strong evidence of the intrinsic emergent properties in human groups, which tend to often outperform their isolated peers of individuals. In addition to the above, many studies looked into the irrational biases in decision making (e.g., cognitive overload and herding in decision making), and bounded rationality (Ariely, 2010; Simon, 1982), along with their negative and positive effects and impact on our daily lives and relationships and the complex collective action. For instance, herd behaviour (Banerjee, 1992), a phenomenon that principally refers to animal group behaviours such as wolf packs, bird flocks and schools of fish, typically when fleeing from a predator, is also observed in human animals as for example, in demonstrations, general strikes, sport events, etc. This seemingly complex (collective) behaviour is in fact simple in nature and emerges without centralised coordination. Yet, it has a large impact on the level of economy (stock markets), the crowd behaviour and its psychology (violence and racial groups) and the day-to-day decision-making, activities and judgement. This human herding behaviour was discussed in the context of the *bandwagon effect* and the *information cascade* (Bikhchandani et al., 1992) principles. Moreover, Gigerenzer studied social intelligence by looking at

how the mind copes with its environment (Gigerenzer, 2000) and, using models of bounded rationality (Gigerenzer and Selten, 2002; Gigerenzer and Goldstein, 1996), showed that cognitive mechanisms capable of successful performance in the real world do not need to satisfy the classical norms of rational inference—and that many models in economics, cognitive science and biology ignore the fact that humans and animals make inferences about the world under limited time and knowledge.

Moving on to political sciences, some studies focused on identifying and analysing factors behind the superiority of collective decision making and its repercussions on the political system (Landemore, 2013), and the social norms (Ostrom, 2014). Collective intelligence was further exploited from a social point of view, in the modelling of collective adaptive systems (Brown and Lauder, 2001). Gruber (Gruber, 2008), for example, considered the *social web* as an ecosystem of participation, and looked at the consequences of combining ideas from the *social web* and *semantic web* by the aggregation of many individual user contributions, showing that this is a crucial step which will open new doors towards the emergence of intelligence in these systems -yet another form of harnessing the wisdom of the crowds.

Castelfranchi (1998) has studied the relationship between sociality (cooperation, competition, groups and organisation) and individual social action and mind. The author highlights the importance of *the social character of the individual action* in reaching an action at collective level. The beliefs, desires and intentions of cognitive agents are shown to be essential for coordination and cooperation along with some emergent pre-cognitive structures and constraints and other emergent forms of cooperation that are required for planning and deliberative agents.

In 2010, Salminen (2012), in his survey on human collective intelligence, classified the relevant studies into three levels of abstraction: the micro-level (collective intelligence as a combination of cognitive and behavioural elements), the macro-level (as a statistical phenomenon witnessed by the crowd) and the level of emergence (as the middle layer translating from micro to macro scale measured using the theories of complex adaptive systems). Engel et al. (2015), and earlier Woolley et al. (2010), empirically showed using a series of online intelligence tests that human groups, just like individuals, have a certain level of intelligence. In their works on computer-mediated collaboration, they revealed that a collective intelligence factor can emerge in human groups and it is strongly correlated with the average social sensitivity of group members, yet not highly dependent on the average or maximum individual intelligence of group members.

As a matter of fact, I have collected a large number of papers and studies which have demonstrated and given empirical evidence of collectives that can outperform non-interactive individuals under different circumstances and settings. In a nutshell, these studies tell us that the performance of human collectives is controlled or influenced by one or more of the group properties/characteristics listed in Table 2.3.

The understanding of some of the underlying characteristics of collective intelligence has significantly enhanced business and management operations. Internet and web-based

<i>Group property</i>	<i>Discussed in</i>
The aggregation details of information collected from the group individuals	(Pentland, 2007; Kasparov and King, 2000; Krause et al., 2011)
The diversity and gender of the group members	(Bonabeau and Meyer, 2001; Hong and Page, 2004; Bonabeau, 2009; Krause et al., 2011; Aggarwal and Woolley, 2013; Woolley et al., 2010)
The members' (virtual) social community formation	(Brabham, 2010; Cachia et al., 2007; Pentland, 2006, 2007)
The network and interaction structure of the group	(Watts and Strogatz, 1998; March, 1991; Watts, 2004; Easley and Kleinberg, 2010; Mason and Watts, 2012; Anicich et al., 2015)
The (collective) decision-making technique used	(Pentland, 2006; Yu et al., 2010; Krause et al., 2011; Charness and Sutter, 2012; Chmait, Dowe, Li, Green and Insa-Cabrera, 2016)
The influence (trust, confidence and bias) between group members and other parties	(Bosse et al., 2006; Cachia et al., 2007; Krause et al., 2011; Koriat, 2012; Satopää et al., 2014)
The motivation and commitment of the members	(Schelling, 2006; Malone et al., 2010; Brabham, 2010)
The communication mode/medium and overall strategy	(Bosse et al., 2006; Anicich et al., 2015; Chmait, Dowe, Green and Li, 2015)
The group's social sensitivity or perceptiveness	(Goleman, 2007; Woolley et al., 2010; Woolley and Bell, 2011; Aggarwal and Woolley, 2013)

Table 2.3: Factors and characteristics impacting the performance of human collectives.

tools have transformed modern enterprise systems by allowing for personalised user content and new forms of user interactions such as collaborative filtering (Nagalakshmi and Joglekar, 2011). Bigger and new types of collaborative projects are made possible, within and between organisations, as working environments become more flexible, interconnected and customisable to project team members. This has increased the competitive advantage of many online businesses like Facebook and Amazon who invest in the power of groups and the crowd. It has been shown (Woolley et al., 2010) through experiments that the team's ability to cooperate and interact effectively is more important for building intelligent groups than the members' individual abilities. Many successful businesses and corporations were aware of this idea even before it was put to test. This was reflected by the different divisional and organisational structures that were implemented by these businesses as for example, flat, hierarchical, product, geographical and matrix organisational structures (Daft, 2012; Tran and Tian, 2013).

In addition to the above, the Internet of Things (IoT) (Gubbi et al., 2013) promoted the emergence of collective intelligence. Nowadays IoT provides highly efficient solutions to many business and societal problems by helping us make more informed and intelligent decisions. The IoT generates massive amounts of information regarding various aspects of our lives. Such information is usually stored in *clouds* and then processed to acquire knowledge used for making intelligent decisions of high impact (Gubbi et al., 2013). The *Inter-cooperative Collective Intelligence* (Xhafa and Bessis, 2014), *Embedded Intelligence*

(Guo et al., 2011) and *Social Web of Intelligent Things* (Console et al., 2011) define example outputs, created by IoT, that had strong impact on the emergence of collective intelligence (Mačiulienė, 2014).

We can easily see how broad is the range of areas and disciplines in which collective intelligence has been studied. There is a long list of terms and concepts which fall under the umbrella of collective intelligence, many which have appeared in this thesis so far, such as: *collaborative systems, community systems, collaborative filtering, crowd-sourcing, crowd-mining, human computation, mass collaboration, prediction markets, recommendation systems, smart mobs, social computing, swarm intelligence, user-powered systems, wisdom of the crowds, social bookmarking, IoT, complex adaptive systems*, etc.

As always, it is desirable to be able to connect these concepts and applications all together, and further find a general way to express, measure and compare the (collective) intelligence of different entities with one another.

2.3 Measuring Individual and Collective Intelligence

Measurement is a fundamental tool for the understanding of (collective) intelligence and its analysis. As we will see in next, the evaluation of intelligence has been a fertile area since the last century. In following sections, I survey many of the important measurement paradigms and experiments that have been used, in different research area and disciplines, to evaluate animal, machine and human cognitive systems.

2.3.1 Evaluating non-human animals

The intelligence of non-human animals has been broadly studied since the early 1900's. This area of research was particularly motivated by Edward Thorndike's experimental studies (Thorndike, 1965; Thorndike et al., 1926) on animal cognition and learning. Numerous apparatus have been used (or adapted) to evaluate the intelligence of animals such as those based on the works of (Tomasello and Call, 1997; Wasserman and Zentall, 2006; Gardner and Gardner, 1969; Bird and Emery, 2009; Hanus et al., 2011; Shettleworth, 2010), just to name a few. Considering the vast number of these evaluation methods, I will only selectively describe a few *intelligence* tests which have been practically utilised for evaluating diverse entities in the animal kingdom. Illustrations of these tests can be found in Figure 2.2.

The first test appearing in Figure 2.2a consists of a binary-tree maze structure (of a laboratory arena) used to measure the rate of information-transmission, and compression, in ants foraging for food (Ryabko and Reznikova, 2009). The ants' *mean duration of transmission of information* on the way to the trough was recorded for different maze structure complexities and food locations. Sample results revealed that ants are able to grasp regularities and to use them for compressing information and transferring it to each other. Likewise, Figure 2.2b depicts a test setup that investigates the ability of honeybees to learn mazes. Results from such experiments showed that honeybee "performance in the

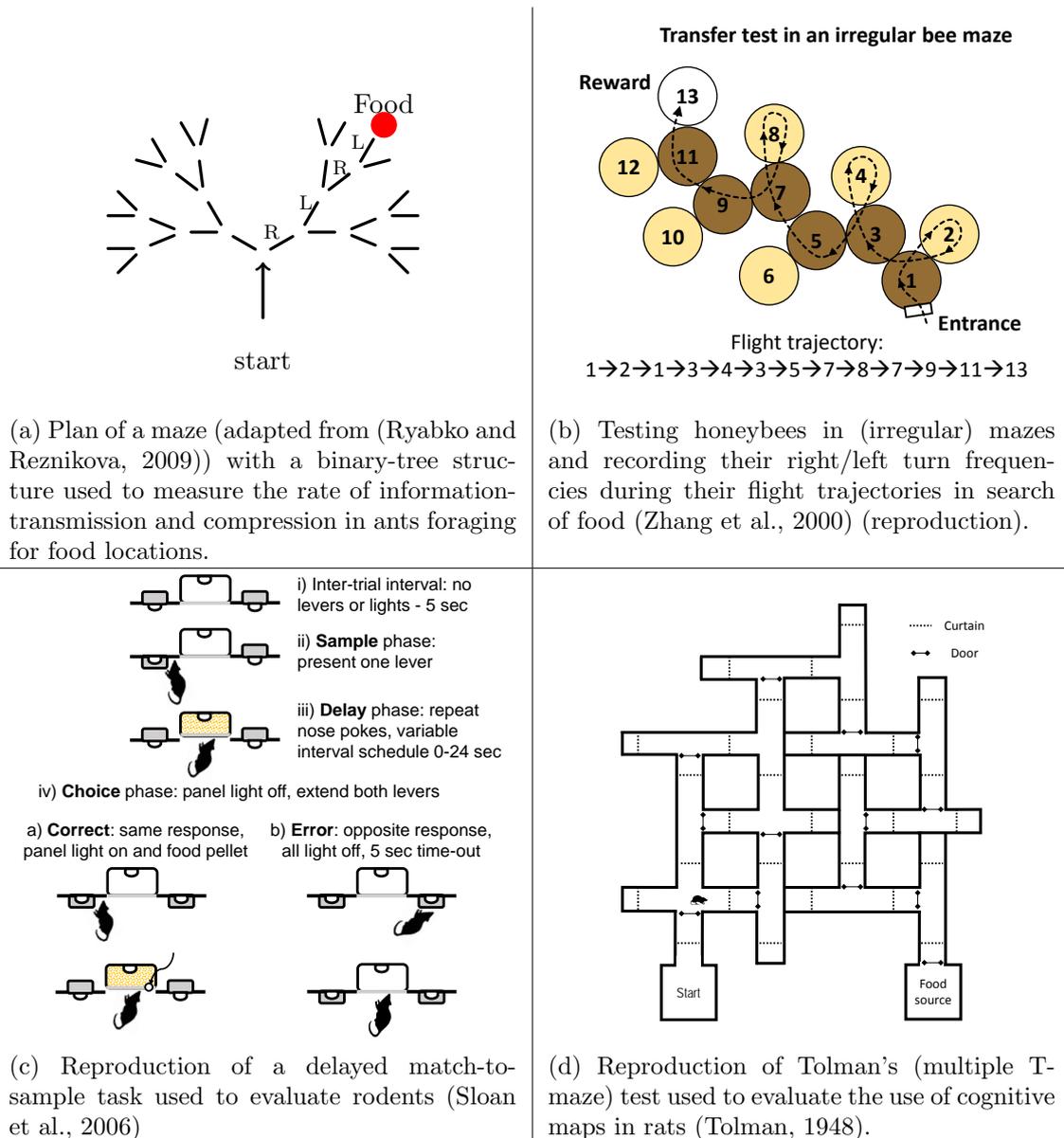


Figure 2.2: A collection of *intelligence* tests used for the evaluation of non-human animals and insects.

various configurations depends on the existence of regularity in the structure of the maze and on the ease with which this regularity is recognised and learned” (Zhang et al., 2000).

Figure 2.2c illustrates a delayed match-to-sample task. This kind of task is used to evaluate the performance (of simultaneous matching, latency, and short term visual memory) of rats (Sloan et al., 2006). A rat is presented with one of two levers in the sample phase, followed by a delay period in which they are retracted, and the rat must choose afterwards the lever that was presented in the sample phase (or in some cases the non-matching one) to receive a food reward. In such experiments the evaluatee could be tested over different visual patterns of varying complexities.

The last picture appearing in Figure 2.2d is a reproduction of Tolman's multiple T-maze test (Tolman, 1948). Experiments over Tolman's test have shown that while rats (and also pigeons) are exploring a maze they develop a cognitive map of its layout. Thus, they

can represent and learn spatial patterns while searching the environment for a particular object. Different maze designs were also used to test for novelty and memory in rats.

It is obvious that the testing environments in Figure 2.2 are different in structure and are highly adapted to the type of the evaluated entity. Nonetheless, broadly speaking, the tasks to be performed are similar to one another, in the sense that they consist of searching and foraging environments, as well as learning and inferring patterns of different complexities. Yet, it is not known how the complexity of one environment compares to another, and how to test for the relative complexity. Apart from environment complexity, another complication in the measurement of intelligence is that not all testing environments are applicable (or could be adapted) to evaluate collectives. While some tests (e.g., Figures 2.2a and 2.2d) might be adequate to evaluate animal collectives, others (e.g., Figures 2.2b and 2.2c) are simply inappropriate, and often we have no *experimental access* to administer such tests in the collective scenario in the first place.

Another *non-test measure* used to estimate the intelligence (or cognition) in animals is the *Encephalization Quotient* (EQ) (Jerison, 2012; Roth and Dicke, 2005). Originally proposed by Jerison (Jerison, 1975, 2012) this quotient reflects the extent to which the brain size of the species to be evaluated deviates from the expected brain size of a standard species in the same taxonomic group. While the EQ has successfully shown that mammals with larger brains are more intelligent in general (Gibson et al., 2001), it has received a few criticisms (Deaner et al., 2007), and it does not apply to some mammals and many groups of invertebrates (Roth and Dicke, 2005, Neural correlates of intelligence).

2.3.2 Evaluating machines and artificial agents

Moving onto the artificial world, perhaps the best introduction into the history of machine intelligence is the imitation game (Turing, 1950) proposed by Turing in the 1950s. While it was once considered as an adequate intelligence test for machines, the test has a few limitations (Oppy and Dowe, 2011) and it is mainly a test of humanness. Afterwards, in the 1960s, an IQ test approach similar to the one used to evaluate humans was repeatedly used to evaluate AI agents. For example, programs that could solve geometric analogy tasks (Evans, 1964), as well as others (Simon and Kotovsky, 1963) able to solve Thurstone’s letter series completion problems (Thurstone and Thurstone, 1941), or even more general IQ tests (Sanghi and Dowe, 2003), were devised. For thorough background on the different computer models solving intelligence test problems, and detailed history of evaluation of (machine) intelligence, see (Hernández-Orallo et al., 2016, Table 4, and Figure 3) and (Hernández-Orallo, 2017) respectively.

Apart from using an IQ test approach, just the same way we can evaluate animals and insects in their natural environment, we can also evaluate (bio- and nature-inspired) computer models and AI agents in artificial environments. For instance, setups similar to those found in Figure 2.2 could be implemented and used to evaluate various types of agent and heuristics. Figure 2.3 highlights exactly this point by illustrating a few simulations of AI agents and collectives operating in artificial environments. Consequently, different ant

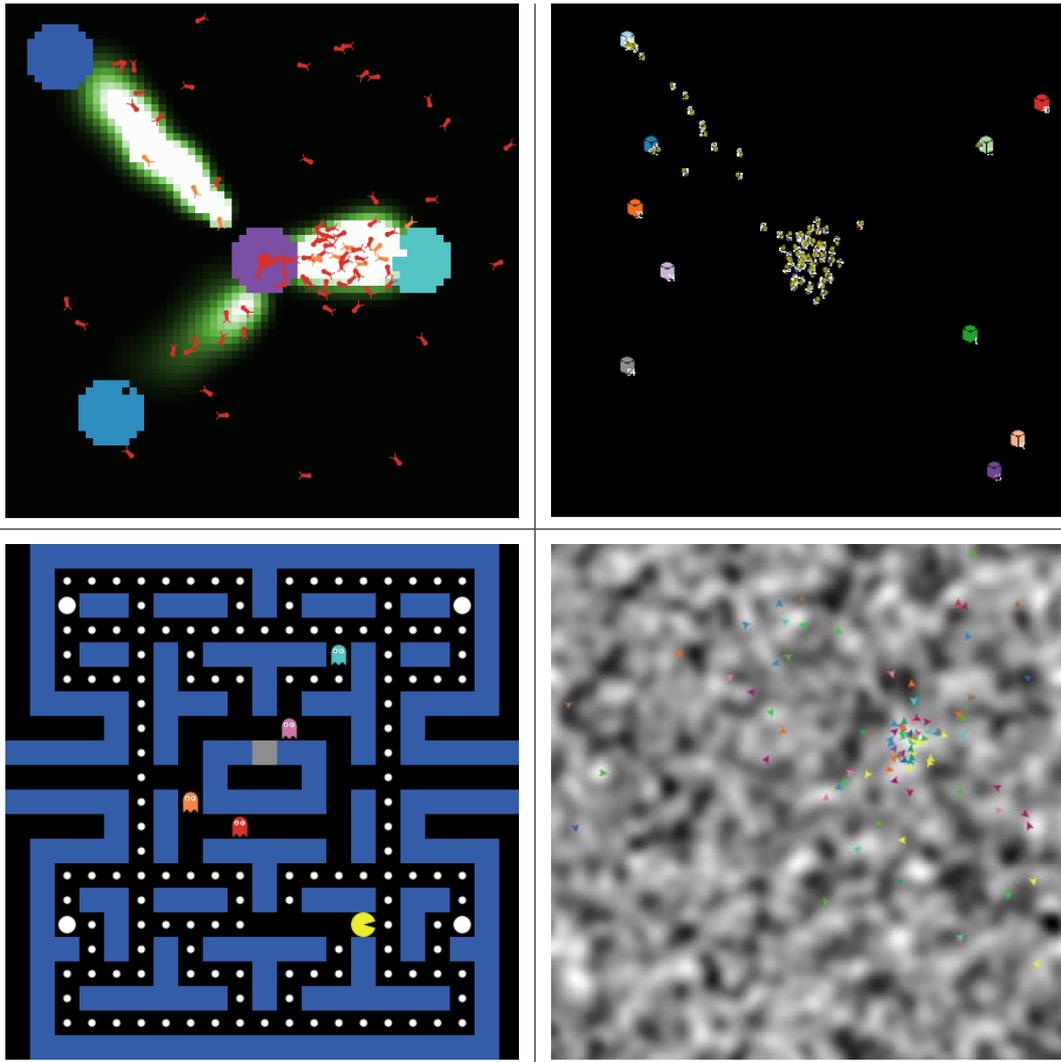


Figure 2.3: Snapshots of exemplar simulations of AI agents and collectives operating in 2D environments taken from (Wilensky, 1999). The top-left image depicts a simulation of an ant-colony (algorithm) foraging for food; the top-right image is a simulation of the intelligent swarm operation of honeybees during their hive-finding process; the bottom-left image depict a simulation of a Ms. Pac-man AI navigating through a maze eating *pellets* and avoiding *ghosts*; the bottom-right image is a snapshot of a particle-swarm-optimisation algorithm searching a fitness landscape trying to optimise a fitness function.

and bee-colony algorithms and heuristics, as well as other general purpose (bio-inspired) algorithms like Particle-Swarm Optimisation (PSO), could be modelled and evaluated over mazes, environments, and tasks of different difficulties and complexities. Taking into consideration the functional similarities between animal swarms and artificial swarms inspired by them, the question of whether evaluating artificial agents is adequate and sufficient for understanding animal cognition and intelligence is still not fully understood. However, at the very least, computer models and simulations have been very informative and provided ingenious explanations of many (human and non-human) animal behaviours and their intelligence characteristics.

Recently several computer games (and game-based competitions) were adopted to evaluate AI agents like Ms. Pac-Man (Lucas, 2007), Unreal Tournament (Hingston, 2010), and the famous Super-Mario Bros competitions (Togelius et al., 2013). Add to that Microsoft’s latest Malmo platform (Johnson et al., 2016) (for running the interactive game Minecraft (Aluru et al., 2015)) proposed as a new experimental world for testing AI agents.

Many of the AI testbeds were influenced by Solomonoff’s general theory of inductive inference (Solomonoff, 1960, 1964a,b), and (Legg and Hutter, 2007)’s definition of *Universal (machine) Intelligence* and the AIXI (Hutter, 2004) agent, dubbed the universal optimal learning agent. For instance, Artificial General Intelligence (AGI) developed as a new branch of research inside artificial intelligence. Interest in (evaluating) general purpose AI (particularly using a reinforcement learning approach) became more popular. In 2010, a set of formal guidelines for designing a universal (machine) intelligence test were presented (Hernández-Orallo and Dowe, 2010). AGI seems to be a promising research area that has engaged many researcher from Facebook (Synnaeve et al., 2016), DeepMind (Beattie et al., 2016) and others (Rosa et al., 2016; OpenAI, 2016; Kempka et al., 2016) in the year 2016. As discussed in (Hernández-Orallo et al., 2017), most of the new AI benchmarks are shifting towards testing *general problem solving ability*. I give below a short overview of these benchmarks following (Hernández-Orallo et al., 2017):

- The **Arcade Learning Environment** is a platform for developing and evaluating general artificial agents using a variety of *Atari 2600* games. This environment allows one to compare reinforcement learning approaches (e.g., (Mnih et al., 2015)), model learning, model-based planning, imitation learning and transfer learning. The small number of games in the Arcade Learning Environment can lead to overspecialisation but this can be overcome using the Video Game Definition Language (VGDL) in which new arcade games can be generated using a malleable set of rules.
- **OpenAI Gym** (Brockman et al., 2016) uses an open-source interface to monitor and compare AI agents interacting with a collection of reinforcement learning tasks. The assessment tasks range from classic control and toy text problems to more advanced algorithmic problems and three-dimensional robots, as well as computer games. Moreover, **OpenAI Universe** (OpenAI, 2016) is a recent software platform that can be used to train and evaluate the performance of AI systems on tasks similar to the ones a human can complete with a computer, and in a similar way to how humans do by looking at screen pixels and operating a (virtual) keyboard and mouse. More than 1000 environments are now ready for evaluating reinforcement learning agent as part of the latest release.
- Microsoft’s **Project Malmo** (Johnson et al., 2016) is a platform in which users have complete freedom to build complex 3D environments within the block-based world of the *Minecraft* video game. This means that the project can support a wide range of experimentation scenarios for evaluating AI agents and a flexible setting for future general AI research. Examples of assessment tasks in this platform include navigation and survival in complex environments and collaborative problem solving.

- GoodAI’s **Brain simulator and school** is a different sort of platform used for the simulation of artificial brain architectures by using existing AI modules like image recognition, working memory and many others.
- **DeepMind Lab** is a platform in which AI agents operate in 3D game-like environments using a first-person viewpoint. The agents can be evaluated over a different range of tasks such as maze navigation and laser-tag games. In a similar approach, the **ViZDoom** (Kempka et al., 2016) research platform allows to monitor reinforcement learning agents interacting in shooting-game scenarios using only the screen buffer.
- Facebook’s **TorchCraft** (Synnaeve et al., 2016) is a library for conducting machine learning research on real-time and real-world strategy games using a high-dimensional action space. This library provides advanced evaluation techniques to reinforcement learning agents by making rewarding actions and planning hierarchically connected and requiring some sort of coordination.
- Facebook’s **CommAI-env** (Mikolov et al., 2015) is another platform for training and evaluating artificial intelligence, this time in the goal of interacting with humans via language. Agent are presented with a range of counting and memory problems of incremental difficulty via a communication-based setup using a bit-level interface.

This new movement towards the evaluation of general AI agents has even given birth to specialised workshops targeting this interesting topic such as the *Evaluating General-Purpose AI (EGPAI)* workshop series (Chmait, Hernández-Orallo, Martínez-Plumed, Stranegård and Thórisson, 2017; Dimitrakakis et al., 2016).

In fact, earlier ideas from (Goertzel and Bugaj, 2009; Adams et al., 2012) proposed a class of environments for teaching and evaluating artificial systems through a process similar to human developmental psychology. The framework for assessing these systems, known as AGI preschool, consists of environments “modelled loosely on preschools used for teaching human children and intended specifically for early-stage systems aimed at approaching human-level” intelligence (Goertzel and Bugaj, 2009)⁴. With regards to the evaluation of probabilistic (Bayesian) estimators, methods such as LNPPP (Dowe, 2008a, Sec. 0.2.7)(Dowe, 2011, Sec. 5.3, p. 936)(Dowe, 2013, Sec. 4.7, p. 24) were proposed to give universal distributions over (environments of) statistical and machine learning problems to compare the efficacy of rival estimators (e.g., AIC vs BIC).

2.3.3 Evaluating human intelligence

⁴One of the challenging competency areas associated with achieving a human-like intelligence is communication (Adams et al., 2012, Tab. 1). Teaching communication (and not simply programming it into an agent) is a difficult task. Chris S. Wallace’s thoughts on trying to communicate with an alien intelligence are discussed in (Dowe, 2008a, Sec. 0.2.5, p. 542, col. 2). Perhaps see also (Dowe, 2013, Sec. 4.4).

The largest number of conclusions on human intelligence in the literature originates from the field of psychometrics, following over 100 years of consideration (Binet and Simon, 1904). To categorise the numerous approaches that were applied in psychometrics to the evaluation of human intelligence, I here distinguish between tests designed to evaluate human personality such as (Goldberg et al., 2006; John et al., 2008), against others which evaluate cognitive abilities. In this chapter I will only look at the latter for two main reasons. Firstly, many of the available personality tests were not validated for reliability (Rust and Golombok, 2014), and secondly, we have no way to quantify the difficulty or complexity of a personality test.

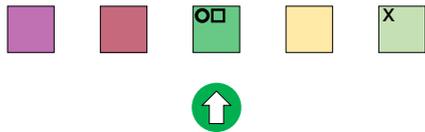


Figure 2.4: A snapshot of the intelligence test (Insa-Cabrera, Dowe, España-Cubillo, Hernández-Lloreda and Hernández-Orallo, 2011) interface used for the evaluation of humans, taken from <http://users.dsic.upv.es/proy/anynt/human1/test.html>.

- $k = 9$: a, d, g, j, ...
- $k = 12$: a, a, z, c, y, e, x, ...
- $k = 14$: c, a, b, d, b, c, c, e, c, d, ...

Figure 2.5: Sample tasks taken from the C-test (Hernández-Orallo, 2000) consisting of completing (inferring the next letter in) alphabetical series and sequences of different complexities 9, 12, and 14.

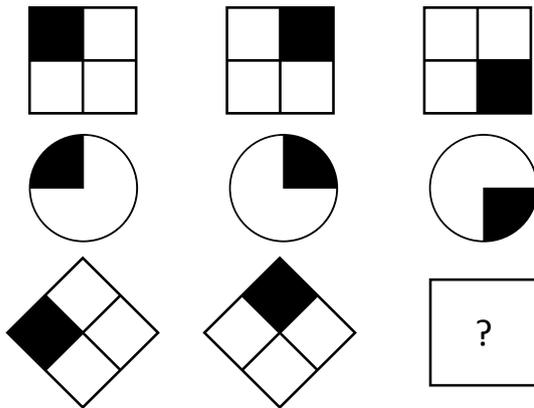


Figure 2.6: A reproduction of a sample intelligence test from Raven's Progressive Matrices test (Raven and Court, 1998). The test evaluates the subject's abstract reasoning abilities and is used as estimate of fluid intelligence.

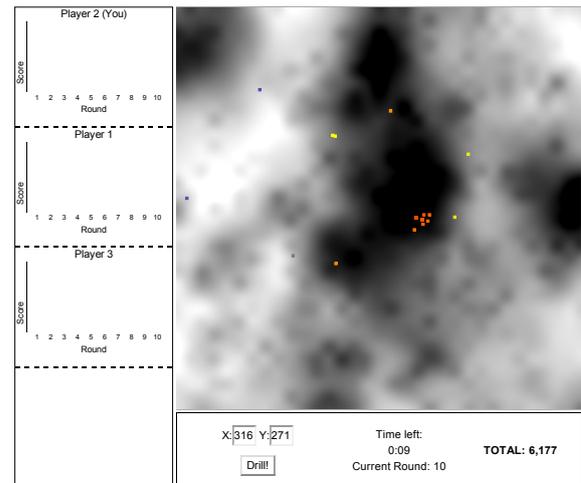


Figure 2.7: Snapshot from a multiplayer online game, called *Wildcat Wells* (<https://volunteerscience.com/test/>), consisting of 2-dimensional desert map with the purpose of searching for oil fields.

Some of the earliest contributions to the area of human intelligence testing were Thurstone's letter series completion problems (Thurstone, 1938; Thurstone and Thurstone, 1941), and later Raven's Progressive Matrices test (Raven and Court, 1998) (Figure 2.6). Both recorded strong correlation with Spearman's general intelligence factor (Spearman, 1904), also known as the "g-factor". Other, more general, tests including a variety of evaluation tasks were also developed in the last century. These tests came to be known as *Intelligence Quotient* or IQ tests. The Stanford-Binet test (Roid, 2003) and the Wechsler

intelligence scales for adults and children (Wechsler, 2008) are good examples of these tests. In order to give a broad picture of the variety of intelligence tests that were administered to humans, I have depicted in Figures 2.4, 2.5, 2.6 and 2.7 four different interfaces (out of many other examples) that these tests can take. Figure 2.4 is a snapshot taken from the Anytime Intelligence Test (Insa-Cabrera, Dowe, España-Cubillo, Hernández-Lloreda and Hernández-Orallo, 2011) (anYnt) interface used for evaluating humans over a series of interactive exercises. The subject, represented by a black circle, should navigate between cells in order to maximise his/her rewards by learning and inferring the cell pattern controlling the allocation of these rewards. The subject is presented with several exercises of different measurable complexities.

Another type of test is depicted in Figure 2.5. In a somewhat similar ambiance and structure to Thurstone’s letter series completion problems, the C-test (Hernández-Orallo, 2000) was designed to measure the ability of comprehension by evaluating a subject over alphabet sequences of measurable complexities. The subject must infer the next letter completing the sequence in order to receive a positive reward. Like the previous test, each of the tasks (sequences) in this test have different and measurable complexities. In fact, some works (Dowe and Hajek, 1997a,b, 1998) preceding the C-test have discussed similar measures of intelligence based on the notion of compression of data and the ability to perform inductive reasoning.

A reproduction of one of Raven’s Progressive Matrices is depicted in Figure 2.6. This test evaluates the subject’s abstract reasoning abilities and is used to estimate her fluid (or ‘g’) intelligence. The main task in this test is to correctly identify the missing element that completes a (visual) pattern. Figure 2.7 is yet another form of test which was used to evaluate humans. The figure shows a snapshot from a multiplayer online game consisting of a two-dimensional desert map in which groups of online players search for (high concentration) oil fields. This game was used to evaluate how the network structure of human groups affects their collaborative learning and exploration performance (Mason and Watts, 2012). Many other tests were devised and used to test human abilities and measure their intelligence, as for example Guilford’s 1967 collection of intelligence tests (Guilford, 1967), the Differential Aptitude Test developed by Bennett (Bennett et al., 1956) and the series of tests developed by the American Council on education psychological examination (Service et al., 1952). However, since many of these tests were originally designed to evaluate individuals, is it valid and reliable to administer these tests to human groups and evaluate their performance?

2.3.4 Individual versus collective intelligence

A fundamental question at the current stage of research in collective intelligence is *where exactly are we in the measurement of collective intelligence?* It is widely accepted that social (contexts for) intelligence cannot be overlooked when evaluating animal troops, swarms, and human groups. Yet many of the tasks and activities designed to evaluate individuals have also been used to evaluate groups. Examples of such tasks are found in (Laughlin, 2011, Chapters 3, 5 and 6).

To highlight the importance of this point, I compare the results from two, almost identical, experiments that were used to measure what was called the *Crowd IQ*. In both experiments, a human-group of 138 members was evaluated over a series of Raven's Progressive Matrices, similar to the one depicted in Figure 2.6, and one answer was elected after aggregating the member's guesses together by majority voting.

The main difference between the two experiments is that the first (Bachrach et al., 2012) was run using members operating in a lab, whereas the second experiment (Kosinski et al., 2012) recruited (crowd-sourced) members using Amazon's Mechanical Turk and only those members with high confidence were encouraged to issue an answer while the others were discouraged. The results from these experiments were significantly different showing higher scores in the second experiment across (and between) the different group sizes. This clearly shows that a test designed to evaluate (a particular type of) individuals might simply not be suitable to evaluate (other types of individuals or) collectives. In case the test was suitable for group evaluation, one must be very careful in interpreting the results and drawing any conclusions or correlations from them. We observe a similar situation in the animal kingdom. For instance, while it might be appropriate to evaluate a honeybee swarm over hive-finding problems (e.g., Figure 2.3, top-right), evaluating more than one honeybee simultaneously in maze environments (e.g., Figure 2.2b) only makes the subjects more confused and leads to an unpredictable random behaviour. Having that in mind, I will look at some of the testbeds that were used to evaluate groups and collectives in various cognitive system.

One of the recent test batteries for collective intelligence⁵ was developed by Engel (2015) following (Woolley et al., 2010). The tests were adapted into an online testing interface where all the evaluated group members could see the input of one another, interact, edit and coordinate their activities in real-time. A few of the tasks used in this online test battery consisted of measuring the accuracy in detecting details and patterns in a large set of data (images, words, etc.), and collectively remembering information from a video-stream, as well as solving matrix puzzle problems. In fact, similar types of exercises on recognition memory were earlier discussed by (Laughlin, 2011, Chapters 3). An example of such exercises consists of having members of a group cooperate together and then answer "True or False" questions after watching a videotape.

I have also presented, earlier above in Figure 2.7, one group intelligence test devised by Mason et al. (2012) that was administered to humans. The authors conducted their test in the form of an online game in which the players must explore a desert map in the search of buried oil-fields. The players were first distributed into various network structures and then interact and collaborate, seeking to maximise their rewards.

Other types of tests take into consideration the age of the evaluated members as, for example, the "New Tanaka B Intelligence Scale" (Tanaka et al., 2003). In this test

⁵Note that I am only interested in empirical tests and evaluation techniques. Other interesting approaches based on the theory of collective action (Ostrom, 1998, 2014) are used to predict the behaviour of groups in different cognitive systems. This is frequently witnessed in the field of game-theory in which formal models are used to evaluate and predict the behaviour (payoff) of competing or cooperative (intelligent rational) agents.

members are evaluated on solving mazes, calculating cubes, replacing figures and numbers, discriminating character strings, completing number series, erasing and completing figures (Uno et al., 2014).

Monitoring the collective behaviour of social systems does not always have to be conducted by means of traditional group intelligence tests or group performance activities. For instance, the behavioural data of various social systems has been monitored using sociometric badges (Olguín and Pentland, 2007), which are personal devices that collect individual behavioural data including audio, location, and movement. Data from these badges were used in the research on influence models. For instance, group interactions and interpersonal influences were recorded and analysed in order to make better predictions about influence within social systems (Pan et al., 2012) and to assess the group performance (Olguín and Pentland, 2010).

Animal groups and swarms are also evaluated as one unit. For instance, we have seen earlier (in Figure 2.2a) a diagrammatic illustration of a laboratory arena binary-maze structure (adapted from (Ryabko and Reznikova, 2009)) which was used to evaluate the performance and ability of ants in transmitting information between each other and compressing patterns while foraging for food. A recent study (Kao et al., 2014) has highlighted the importance of collective evaluation of animals by empirically measuring and comparing the learned performance of animals within groups as opposed to individual evaluation. This has led to new results showing how associative learning functions within a social context in animals.

Independently, other types of experiments were designed to test for observational learning in animals. In many cases, animals were not evaluated in the same (physical) environment but instead in a similar world making sure the testee can observe other members/demonstrators. For instance, it was shown (Fiorito and Scotto, 1992) that untrained *Octopus Vulgaris* learn remarkably by observation, irrespective of the object chosen by the demonstrators as the reward. Similarly, other types of animals and cognitive systems such as humans and robots can learn by imitation and social learning (Nehaniv and Dautenhahn, 2007). For a discussion on the history, definition, and interpretation on experimental data on imitation in animals see (Galef Jr, 1988). Specialised tests have been conducted to measure coordination between partners in elephants, as opposed to imitation. The tests revealed that elephants not only cooperate with their partner, but also understand the logic behind teamwork (Plotnik et al., 2011).

Moreover, non-test approaches to the measurement of intelligence have been relatively useful to understand social intelligence in animals. By interpreting the Encephalization Quotient (Jerison, 2012; Roth and Dicke, 2005), interesting connections were discovered between the social complexity and brain size of a species. For instance, it was shown (Shultz and Dunbar, 2006) that the relative brain size was independently associated with sociality and social complexity of animal groups, and also related to the nature of social relationships as well as the total number of individuals in a group.

Again in the case of artificial agents, robots and multiagent systems in general, many (social) games featuring more than one agent were used for evaluation. For instance,

games like RoboCup (Kitano et al., 1997) and other predator-prey games were used to evaluate the performance of multiagent systems, as well as their collective behaviour—e.g., investigating the role of co-evolution in the context of evolutionary robotics (Nolfi and Floreano, 1998). Moreover, Artificial Neural Networks can be used to emulate and predict the behaviour of a large number of interacting (artificial/physical) agents (neurons). In this scenario, a population of agents can be modelled as a dynamic/evolving neural network (Byrski and Kisiel-Dorohinicki, 2003) that exhibits some kind of intelligence at the level of the collective after different periods of interaction steps.

What is more, hybrid collectives were also evaluated. By hybrid I mean collectives containing members from more than one cognitive system. For instance some experiments were performed on a robotic rat which was evaluated with real rats (Ishii et al., 2006), acting as their social partner. Similar experiments were also performed on schools of (real and artificial) fish (Swain et al., 2012).

At first glance, one tends to ignore the distinction between humans and machines when it comes to evaluation. With the ubiquity of human-computer interaction in our daily activities, and the increase of human access to technology, it has become intuitive to look at humans using technological devices as one entity. For instance, many (e.g., programming and online multiplayer) competitions rely on the skills of the evaluated subjects in using and manipulating technological devices or software systems. While the evaluated subject in such tests is a human, the actual test involves both human and machine entities—blurring the distinction between the two cognitive systems.

Chapter 3

Factors of Collective Intelligence

No one can whistle a symphony. It takes a whole orchestra to play it.

—H. E. Luccock

The research in this chapter has been published in the following articles:

- Nader Chmait, David L. Dowe, Yuan-Fang Li, David G. Green, and Javier Insa-Cabrera (2016). *Factors of collective intelligence: How smart are agent collectives?*, Proceedings of the 22nd European Conference on Artificial Intelligence (ECAI), The Hague, The Netherlands, Vol. 285 of Frontiers in Artificial Intelligence and Applications, IOS press, pp. 542-550.
<http://ebooks.iospress.nl/volumearticle/44798>
- Nader Chmait, Yuan-Fang Li, David L. Dowe, and David G. Green (2016). *A dynamic intelligence test framework for evaluating AI agents*, Proceedings of the 1st International Workshop on Evaluating General-Purpose AI (EGPAI 2016), European Conference on Artificial Intelligence (ECAI) 2016, pp. 1-8. <http://www.ecai2016.org/content/uploads/2016/08/W14-EGPAI-2016.pdf>

3.1 Overview

The dynamics and characteristics behind intelligent cognitive systems lie at the heart of understanding, and devising, successful solutions to a variety of multiagent problems. Despite the vast literature on collective intelligence, much of which I surveyed in the previous chapter, important questions like “how does the effectiveness of a collective compare to that of its isolated members?” and “are there some general rules or properties shaping the spread of intelligence across various cognitive systems and environments?” remain somewhat of a mystery.

In this chapter I develop the idea of collective intelligence by giving some insight into a range of factors hindering and influencing the effectiveness of interactive cognitive systems. I measure the influence of each of the examined factors on intelligence independently, and

empirically show that collective intelligence is a function of them all simultaneously. I further investigate how the organisational or network structure of equally sized groups shapes their effectiveness. The outcome is fundamental to the understanding and prediction of the collective performance of multiagent systems, and for quantifying (the emergence of) intelligence over different environmental settings.

3.2 Introduction

As discussed in the previous chapter, collective intelligence emerges in all sorts of cognitive systems, from natural (e.g., animal and human) to artificial (e.g., software agents and robotics), by cause of diverse social organizations (human societies, efficient markets, social insect colonies, group collaborations via the world wide web, etc.). It seems that the complex structure and operation of these systems hinder our understanding of the dynamics and characteristics behind intelligent collectives, which are fundamental for devising successful models and solutions to a variety of multiagent problems. Despite the broad literature on collective intelligence, the questions “how does the effectiveness of a collective compare to that of its isolated members?” and, more importantly, “are there some general rules shaping the spread of intelligence which can be perceived across different cognitive systems and environments?” remain somewhat of a mystery.

Now imagine we had a series of performance tests which we can administer on different cognitive systems, could we then disclose any patterns or factors at all, explaining the emergence of intelligence among all of these systems?

In this chapter, I will give insight into the main components and characteristics of collective intelligence by applying formal tests for the purpose of measuring (or quantifying) the influence of several factors on the collective behaviour and the accuracy of a group of agents. I also discuss how these results compare to individual agent scenarios. I attempt to uncover some of the dynamics and circumstances behind intelligent collectives in general, hoping this would reinforce the understanding and prediction of the behaviour of groups, by bringing some new results into the AI community.

3.3 On Collective Intelligence

Earlier studies (Engel et al., 2015; Woolley et al., 2010) have revealed that a collective intelligence factor can emerge in human groups. We know that collectives can outperform individuals, and further that their performance is controlled by one or more of:

- (i) their organisational or network structure (Mason and Watts, 2012; Child, 1972; Mintzberg, 1979),
- (ii) the information aggregation details among their individuals (Bettencourt, 2009), and
- (iii) the diversity between their members (Hong and Page, 2004; Hashemi and Endriss, 2014).

Crowd-computing and crowd-sourcing (Poesio et al., 2015; Klein and Garcia, 2015; Bonabeau, 2009) methodologies are excellent examples of collective intelligence that harness the *wisdom of the crowd* (Surowiecki, 2005).

After carefully looking at the literature on collective intelligence, including the surveyed works from the previous chapter, I filter a set of factors or features from these works—that are not coupled to one particular cognitive system, problem or environment—which are intimately relevant to the performance of collectives. Some of these factors are the number of members in a group, the communication or interaction protocol, as well as the difficulty of the environment. Curiously, there are some other factors which are often relatively neglected, such as the reasoning/learning speed of the agents and the interaction time of the collective as a whole. These features, in addition to some hypothetical combinations of them (grouped in ellipses) are depicted in Figure 3.1. It is important to note that there exist other important features that shape the intelligence of groups which are not captured in Figure 3.1. For example, it has been shown (Woolley et al., 2010) that the number of socially competent individuals in a group is linked to the group’s effectiveness in solving difficult problems and that women tend to have better social skills than men on average. While this clearly is an important observation, it does not scale to all sorts of agent cognitive types (e.g., we cannot control for gender in machines) and thus will not be investigated here given that it does not align with the objectives of this chapter.

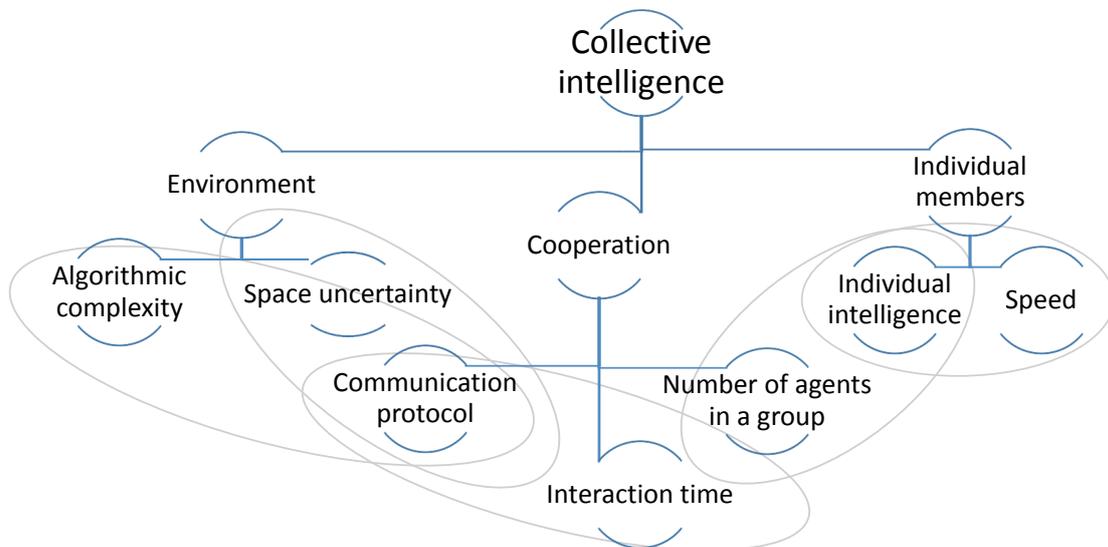


Figure 3.1: Factors and features relevant to the notion of collective intelligence that are perceived throughout various cognitive systems (Chmait, Dowe, Li, Green and Insa-Cabrera, 2016, Fig. 1). Some hypothetical relationships between the factors are also grouped in ellipses.

It is not known *in which circumstances* and *how much* the features/factors appearing in Figure 3.1 individually influences the intelligence of the group, let alone the simultaneous influence of multiple features combined, which is what I attempt to *quantitatively* investigate in this chapter. To answer these questions, we need a dynamic environment in which we can simulate and assess the influence of these factors on the performance of

agents (of different behaviours) operating under different environmental and interactive settings.

Therefore, in the next section I give a brief background on the main approaches that have been used to assess the intelligence of machines or artificial agents and discuss their appropriateness for evaluation. For detailed description of these and other approaches refer back to Section 2.3.2. I continue by introducing my methodology for evaluating intelligence using an agent-environment framework (Section 3.5). I then introduce and describe in Section 3.6 the structure of the Λ^* (Lambda Star) testing environment, which implements the anYnt test (Hernández-Orallo and Dowe, 2010), that will be used to evaluate the effectiveness of agents. I also discuss the details of the test implementation, its setup and parameters (Section 3.7). The agent behaviours to be evaluated, and their communication and interaction protocols, are described in Sections 3.8 and 3.9 respectively. After I present my experimental setup in Section 3.10, I discuss and analyse the results from the conducted experiments (in Sections 3.11 and 3.12) by making a series of observations on how the intelligence of the evaluated agents was influenced by a collection of examined factors. I also draw some conclusions connecting the research outcomes. I then discuss in Section 3.13 some alternative testing environments that might be useful to quantify the performance of artificial agents and raise some arguments and considerations relevant to the robustness of multiagent performance tests. I conclude in Section 3.14 with a brief summary of the outcomes from this chapter.

3.4 Evaluating Artificial Agents: A Short Background

Perhaps a good start to understand the history of machine intelligence would be to take a look back at the imitation game (Turing, 1950) proposed by Turing in the 1950s, where the idea is to have one or more human judges interrogating a program through an interface, and the program is considered intelligent if it is able to fool the judges into thinking that they are interrogating a human being. While this was once regarded as an intelligence test for machines, it has limitations (Oppy and Dowe, 2011) and is mainly a test of humanness. The machine intelligence quotient (MIQ) using fuzzy integrals was presented (Bien et al., 2002) in 2002. However, determining a universal *intelligence quotient* for ranking artificial systems is not very practical and is almost unmeasurable due to the vast non-uniformity in the performances of different types of artificial systems. Several studies, such as (Chaitin, 1982; Dowe and Hajek, 1997a,b, 1998; Dowe et al., 2011; Hernández-Orallo and Minaya-Collado, 1998; Sanghi and Dowe, 2003), investigated the relevance of compression (Dowe and Hajek, 1997a, Sections 2 and 4), pattern recognition, and inductive inference (Dowe and Hajek, 1998, Section 2) to intelligence. Shortly after (Chaitin, 1982; Dowe and Hajek, 1997a,b, 1998; Hernández-Orallo and Minaya-Collado, 1998) came the C-test (Hernández-Orallo, 2000), which was one of the first attempts to design an intelligence test consisting of tasks of quantifiable complexities. However, the test was static (non-dynamic) and it did not fully embrace the vast scope implicit in the notion of intelligence. In 2007, Legg and Hutter proposed a general definition of *universal (machine) intelligence* (Legg and Hutter, 2007), and three years later a test influenced by this new definition, namely the *Anytime*

Universal Intelligence Test, was put forward by Hernandez-Orallo and Dowe (2010) in order to evaluate intelligence. The test was designed to be administered on different types of cognitive systems (human, animal, artificial), and examples environment classes illustrating the features of the test were also suggested in (Hernández-Orallo and Dowe, 2010).

It is clear that many studies have looked into the intelligence of artificial agents. However, beyond single agent intelligence (e.g., (Bien et al., 2002; Hernández-Orallo, 2000; Legg and Hutter, 2007; Hernández-Orallo and Dowe, 2010)), no *formal intelligence tests* were developed in the purpose of *quantifying the intelligence of groups of interactive agents against isolated (non-interactive) agents* across different interactive settings—which is one of the motivations behind this chapter. Yet, before proceeding with the description of this work, I address an important question that might come to a reader’s mind: can’t one simply evaluate and compare artificial systems over any given problem or environment from the literature?

There are several reasons why we can’t do that, most of them were studied and examined in (Hernández-Orallo and Dowe, 2010) and (Legg and Veness, 2013). I give a brief summary of some of these principles. Firstly, there is a risk that the choice of the environment used for evaluation is strongly *biased*, and that it favours particular types of agents while it is unsuitable for others. Furthermore, the environment should handle *any level of intelligence* in the sense that dull or brilliant, and slow or fast, systems can all be adequately evaluated. The test needs to be practical and preferably return a (valid) score after *being stopped at any time-period*, short or long. (Legg and Veness, 2013, Sec. 3.2) also raise some practical issues regarding developing a real-world performance metric showing that one has to somehow control or avert generating non-halting *programs* that are used for defining testing environments. Besides, not every evaluation metric is a formal intelligence test or even at a minimum, a reliable performance metric. For instance, the testing environment should be non-ergodic but reward-sensitive with no sequence of actions leading to heaven or hell scoring situations (e.g., states where, respectively, an agent is always positively rewarded or penalised for all viable actions). The environment should also be *balanced* in the sense that it must return a null reward to agents with a random behaviour.

Although one principal advantage from this work is the measurement of intelligence of (artificial) agents, the outcome also has implications for the design and analysis of agent-based systems. This is because the work provides an opportunity to predict the effectiveness (expected performance) of existing (artificial) systems under different collaboration scenarios and problem complexities. In other words, it is one way of looking at the dynamics behind collective intelligence in multiagent systems.

3.5 Agent-Environment Framework

A common setting in most approaches to measuring intelligence is to evaluate a subject over a series of problems of different complexities and return a quantitative measure or score reflecting the subject’s performance over these problems (Hernández-Orallo

and Dowe, 2010). In artificial systems and simulations, the agent-environment framework (Legg and Hutter, 2007) is an appropriate representation for this matter. For instance, this framework allows us to model and abstract any type of interactions between agents and environments. It also embraces the *embodiment thesis* (Brooks, 1991) by embedding the agents in a flow of observations and events generated by the environment.

Here I define an environment to be the world where an agent π , or some group of agents $\{\pi_1, \pi_2, \pi_3, \dots, \pi_n\}$, can interact using a set of *observations*, *actions* and *rewards*. The environment generates observations from the set of observations \mathcal{O} , and rewards from $\mathcal{R} \subseteq \mathbb{Q}$, and sends them to all the agents. Then, each agent performs actions from a limited set of actions \mathcal{A} in response. An iteration or step i stands for one sequence of *observation-action-reward*. An observation at iteration i is denoted by o_i , while the corresponding action and reward for the same iteration are denoted by a_i and r_i respectively. The string $o_1a_1r_1o_2a_2r_2o_3a_3r_3$ is an example sequence of interactions, over three consecutive iterations, between one agent and its environment. An illustration of the agent-environment

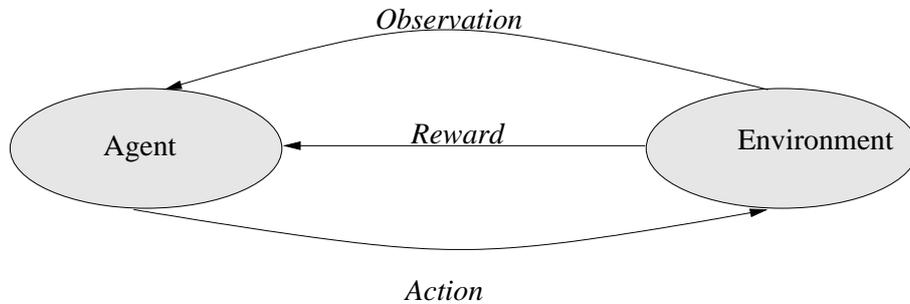


Figure 3.2: Agent-Environment Framework (Legg and Hutter, 2007)

framework is given in Figure 3.2. I define the multiagent-environment framework as an extension of the above such that $o_{i,j}$, $a_{i,j}$ and $r_{i,j}$ are respectively the observation, action and reward for agent π_j at iteration i . The order of interactions is as follows:

1. The environment sends observations to all the agents at the same time.
2. Then, the agents interact (and communicate) and perform corresponding actions.
3. Finally, the environment provides the agents back with rewards.

For instance, the first interaction of agents π_1, π_2 and π_3 in the multiagent-environment setting, denoted by $o_1a_1r_1$, is equivalent to $o_{1,1}o_{1,2}o_{1,3}a_{1,1}a_{1,2}a_{1,3}r_{1,1}r_{1,2}r_{1,3}$.

3.6 Intelligence Test

In order to assess the performances of agents, whether in isolation or collectively, there is a need for an environment over which we can run formal intelligence tests (of measurable complexities) on artificial agents using the recently described framework. Many environments and tasks might be relevant to the evaluation of agent-based systems. However, not every evaluation metric can be used as a formal (general) intelligence test. From the

set of appropriate ones, I have chosen for my purpose (an extension of) the Λ environment class, described in (Insa-Cabrera, Hernández-Orallo, Dowe, España and Hernández-Lloreda, 2012, Sec. 6) and (Hernández-Orallo and Dowe, 2010), which implements the theory behind the “Anytime Universal Intelligence Test” (anYnt) (Hernández-Orallo and Dowe, 2010).

There are many reasons for selecting the *Anytime Universal Intelligence Test* (anYnt) and the Λ environment in particular for evaluation. On one hand, they are derived from a **mathematical background** and rely on well-grounded theoretical foundations. On the other hand, they have been practically used to **evaluate diverse kinds** (including machines) of **entities** (Insa-Cabrera, Benacloch-Ayuso and Hernández-Orallo, 2012; Insa-Cabrera, Dowe, España-Cubillo, Hernández-Lloreda and Hernández-Orallo, 2011). This means that the same type of test can be administered on human, animal and artificial agents using significantly different interfaces while guaranteeing an equivalent measure of difficulty. More importantly, such selection embraces all of the concerns raised in Sections 3.3 and 3.4 regarding the measurement of intelligence. For instance, the Λ environment is an unbiased dynamic setting that can be stopped at any time and used to quantitatively assess the effectiveness of artificial (and other types of) agents.

3.6.1 The Λ^* (Lambda star) environment

In this section I introduce the Λ^* (Lambda Star) environment class which is my extension of the Λ environment, described in (Hernández-Orallo and Dowe, 2010, Sec. 6.3) and (Insa-Cabrera, Hernández-Orallo, Dowe, España and Hernández-Lloreda, 2012), that focuses on a restricted—but important—set of tasks in AI and intelligence measurement in general.

The general idea is to evaluate an agent that can perform a set of actions, by placing it in a grid of cells with two special objects, *Good* (\oplus) and *Evil* (\ominus), travelling in the space using movement patterns of measurable complexities. The rewards are defined as a function of the position of the evaluated agent with respect to the positions of \oplus and \ominus .

Structure of the test

An environment space is defined as an m-by-n grid-world populated with objects from $\Omega = \{\pi_1, \pi_2, \dots, \pi_x, \oplus, \ominus\}$, the finite set of objects. The set of evaluated agents $\Pi \subseteq \Omega$ is $\{\pi_1, \pi_2, \dots, \pi_x\}$. Each element in Ω can have actions from a finite set of actions $\mathcal{A} = \{up-left, up, up-right, left, stay, right, down-left, down, down-right\}$. All objects can share the same cell in the environment space at the same time except for the special objects \oplus and \ominus where, in this case, one of them is randomly chosen to move to the intended cell while the other one keeps its old position. In the context of the agent-environment framework, described in Section 3.5, a test episode consisting of a series of ϑ interactions $o_i a_i r_i$ such that $1 \leq i \leq \vartheta$ is modelled as follows:

1. the environment space is first initialised to an m-by-n toroidal grid-world, and populated with a subset of evaluated agents from $\Pi \subseteq \Omega$, and the two special objects \oplus and \ominus ,

2. the environment sends to each agent a description of its range of 1 *Moore neighbour* cells (Gray, 2003; Weisstein, 2015) and their contents, corresponding to the rewards in these cells, as an observation,
3. the agents (communicate/interact and) respond to the observations by performing an action in \mathcal{A} , and the special objects perform the next action in their movement pattern,
4. the environment then returns a reward to each evaluated agent based on its position (distance) with respect to the locations of the special objects,
5. this process is repeated again from point #2 until a test episode is completed, that is when $i = \vartheta$.

The Λ^* environment consists of a toroidal grid space in the sense that moving off one border makes an agent appear on the opposite one. Consequently, the distance between two agents is calculated using the surpassing rule (toroidal distance) such that, in a 5-by-5 grid space for example, the distance between cell (1, 3) and (5, 3) is equal to 1 cell. An illustration of a sample Λ^* environment is given in Figure 3.3.

Rewarding function

The environment sends a reward to each evaluated agent from the set of rewards $\mathcal{R} \subseteq \mathbb{Q}$ where $-1.0 \leq \mathcal{R} \leq 1.0$. Given an agent π_j , its reward $r_j^i \in \mathcal{R}$ at some test iteration i can be calculated as:

$$r_j^i \leftarrow \frac{1}{d(\pi_j, \oplus) + 1} - \frac{1}{d(\pi_j, \ominus) + 1}$$

where $d(a, b)$ denotes the (toroidal) distance between two objects a and b . Recall that an agent does not have a full representation of the space and only receives observations of its (range of 1 Moore (Gray, 2003; Weisstein, 2015)) neighbourhood. Therefore, I constrain the (positive and negative) rewards an agent receives from the environment (as a function of its position with respect to \oplus and \ominus respectively) as follows:

- The positive reward π_j receives at each iteration is calculated as $1/(d(\pi_j, \oplus) + 1)$ if $d(\pi_j, \oplus) < 2$, or 0 otherwise.
- Likewise, its negative reward at that iterations is $-1/(d(\pi_j, \ominus) + 1)$ if $d(\pi_j, \ominus) < 2$, or 0 otherwise.

The total reward, r_j^i of agent π_j at iteration i , is the sum of its positive and negative rewards received at that iteration.

3.6.2 Algorithmic complexity

The main assessment task in the Λ^* environment is for the evaluatee to find an action-selection policy (or, equivalently, find a sequence of actions) maximising its rewards. This

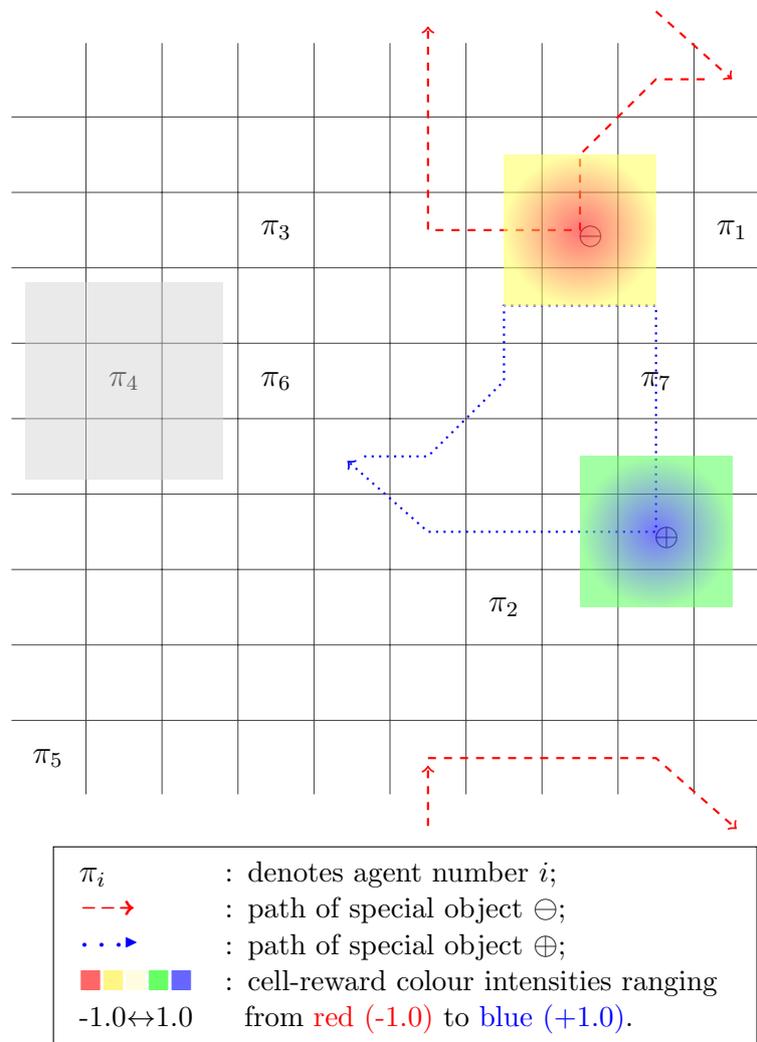


Figure 3.3: A diagrammatic representation of a sample 10-by-10 Λ^* environment, used to implement the theory behind the Anytime Universal Intelligence Test (anYnt) (Hernández-Orallo and Dowe, 2010). Seven agents (π_1, \dots, π_7) are depicted in the toroidal environment space as well as two special objects, \oplus and \ominus , each navigating according to a fixed movement pattern (denoted by the dashed arrow). An agent π_i can only observe or perceive its (range of 1) Moore neighbourhood. For example, agent π_4 can only observe (the cells falling in) the grey shaded area surrounding it (Chmait, Dowe, Li, Green and Insa-Cabrera, 2016, Fig. 2).

entails learning the rewarding pattern generated by the environment and using this knowledge to predict future reward allocations (by intercepting \oplus) while avoiding injury (by escaping \ominus).

Thus, I regard the Kolmogorov complexity (Li and Vitányi, 2008)⁶ of the movement patterns of the special objects as a measure of the algorithmic information-theoretic complexity $K(\mu)$ of the environment μ in which they operate. For instance, a Λ^* environment of high Kolmogorov complexity is sufficiently rich and structured to generate complicated (special object) patterns/sequences of seeming randomness.

The Kolmogorov complexity (Li and Vitányi, 2008) (Definition 1) of a string x is the length of the shortest program p that outputs x over a reference (Turing) machine U .

Definition 1 Kolmogorov Complexity

$$K_U(x) := \min_{p: U(p)=x} l(p)$$

where $l(p)$ denotes the length of p in bits, and $U(p)$ denotes the result of executing p on a Universal Turing machine U .

For brevity, I use the term *algorithmic complexity* to denote the algorithmic information-theoretic complexity $K(\mu)$ of an environment μ . Next, I give an example of how to approximate the value of $K(\mu)$. Assume that the special object \oplus moves in a 5-by-5 grid space. It has an ordered (and repeating) movement pattern travelling between cells with indices: 7, 3, 4, 9 and 8 (corresponding to the greyed out cells appearing in Figure 3.4) such that, in a 25-cell grid, indices 1, 2 and 6 correspond respectively to cells with coordinates (1, 1), (1, 2) and (2, 1) and so on (see Figure 3.4). Also assume that the number of time steps in one test episode ϑ is 20 iterations. Using notions from algorithmic information theory, namely Kolmogorov complexity, I measure the *algorithmic complexity* of the environment $K(\mu)$ in which \oplus operates as the length of the shortest program that outputs the sequence 73498734987349873498 (of length ϑ). Since the Kolmogorov complexity is uncomputable, I measure the Lempel-Ziv complexity (Lempel and Ziv, 1976) of the movement patterns as a practical alternative⁷ (and possibly an approximation) to $K(\mu)$ as suggested in (Lempel and Ziv, 1976; Evans et al., 2002) and (Kaspar and Schuster, 1987, Sec. II).

⁶I repeat here a footnote from one of my earlier joint papers. “The concept of Kolmogorov complexity or algorithmic information theory (AIT) is based on independent work of R. J. Solomonoff (Solomonoff, 1960, 1964a,b) and A. N. Kolmogorov (Kolmogorov, 1965) in the first half of the 1960s, shortly followed by related work by G. J. Chaitin (Chaitin, 1966, 1969). The relationship between this work and the Minimum Message Length (MML) principle (also from the 1960s) (Wallace and Boulton, 1968) is discussed in (Wallace and Dowe, 1999), (Wallace, 2005, Chapter 2 and Section 10.1) and (Dowe, 2011, Sections 2 and 6)” (Chaitin, Li, Dowe and Green, 2016).

⁷ There is no general algorithm that can determine the Kolmogorov complexity of a given string (Solomonoff, 1964a; Kaspar and Schuster, 1987). Nevertheless, Kaspar and Schuster explicitly recognise in the first paragraph of (Kaspar and Schuster, 1987, Sec. II) that (Lempel and Ziv, 1976) provide an appropriate alternative measure of the Kolmogorov complexity of a string by calculating a number $c(n)$, instead of the length of the program which generates a given string of length n , which is a useful measure of this length. Furthermore, (Evans et al., 2002) discuss why the Lempel Ziv 78 Universal compression algorithm (Ziv and Lempel, 1978) is a computationally efficient method towards approaching the *estimation* of the Kolmogorov complexity.

1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20
21	22	23	24	25

Figure 3.4: Encoding the movement pattern of special objects: a conceptual representation of a 5-by-5 grid space with cell indices ranging from 1 to 25. The movement pattern of a special object navigating in this space can be encoded as the list of ordered (cell) indices along which this object travels during a test episode (Chmait, Li, Dowe and Green, 2016, Fig. 3).

Note that, at one test episode, the movement patterns of the special objects \oplus and \ominus are different but (algorithmically) equally complex making sure the rewards are balanced (Hernández-Orallo, 2010). The recurrent segment of the movement pattern is at least of length one and at most $\lfloor \vartheta/2 \rfloor$, cyclically repeated until the final iteration (ϑ) of the test.

3.6.3 Search space complexity

The size of the search space is an important characteristic that is frequently discussed in the context of the difficulty of planning and optimisation problems. This characteristic has prompted the development of different approaches towards solving these problems.

I measure the search space complexity $\mathcal{H}(\mu)$ as the amount of *uncertainty* in environment μ , expressed by Shannon’s entropy (Shannon, 1948). Let N be the set of all possible states of an environment μ such that a state s_μ , is the set holding the current positions of the special objects $\{\oplus, \ominus\}$ in the m-by-n space. Thus the number of states $|N|$ increases with the increase in the space dimensions m and n, and it is equal to the number of permutations ${}^{m \times n}P_2 = \frac{(m \times n)!}{(m \times n - 2)!}$. The entropy is maximal at the beginning of the test as, from an agent’s perspective, there is complete uncertainty about the current state of μ . Therefore $p(s_\mu)$ follows a uniform distribution and is equal to $1/|N|$. Using \log_2 as a base for our calculations, we end up with: $\mathcal{H}(\mu) = - \sum_{s_\mu \in N} p(s_\mu) \log_2 p(s_\mu) = \log_2 |N|$ bits.

3.6.4 Further thoughts on complexity

The *algorithmic* and *search space* complexities of an environment could be combined into a higher level complexity measure of the whole environment. This new measure can be very useful to weight environments that are used for the measurement of (universal) intelligence. Nonetheless, having two separate measures of complexity also means that we can quantify the individual influence of each class (or type) of complexity on the performance of agents. This approach appears to be particularly useful for evaluating the factors influencing the performance of agent collectives as these collectives can exhibit different behaviours in response to changes in the measures of each class of environment complexity.

Note here the importance of measuring the complexity of the testing environment as this will serve, later on, for quantitatively comparing and analysing the impact of environment complexity on the effectiveness of the evaluated agents. This is not the case in the majority of intelligence tests in which different tasks are qualitatively ranked (e.g., as easy or hard)—thus inhibiting our ability of conducting accurate evaluations.

3.6.5 Intelligence score

The metric of (individual agent) universal intelligence defined in (Hernández-Orallo and Dowe, 2010, Definition 10) was extended into a collective intelligence metric (Definition 3) returning an average reward accumulation (per-agent) measure of success (Definition 2) for a group of agents Π , over a selection of Λ^* testing environments (Section 3.6.1).

Definition 2 *Given a Λ^* environment μ and a set of (isolated or interactive) agents $\Pi = \{\pi_1, \pi_2, \dots, \pi_n\}$ to be evaluated, the (average per-agent per-iteration) reward $\tilde{R}_{\Pi, \mu, \vartheta}$ of Π over one test episode of ϑ -iterations is calculated as:*

$$\tilde{R}_{\Pi, \mu, \vartheta} = \frac{\sum_{j=1}^n \sum_{i=1}^{\vartheta} r_j^i}{n \times \vartheta}.$$

Definition 3 *The (collective) intelligence of a set of agents Π is defined as:*

$\Upsilon(\Pi) = \frac{1}{\omega} \sum_{\mu \in L} \tilde{R}_{\Pi, \mu, \vartheta}$, where L is a set of ω environments $\{\mu_1, \mu_2, \dots, \mu_\omega\}$ such that $\forall \mu_t, \mu_q \in L: \mathcal{H}(\mu_t) = \mathcal{H}(\mu_q)$, and $\forall \mu_i \in L, K(\mu_i)$ is extracted from a range of (movement pattern) algorithmic complexities in $]1, K_{\max}]$.

Note the use of the same search space complexity, but different algorithmic complexities, in the intelligence measure defined in Definition 3. As we will see in later sections, the reason behind this is to allow for running controlled experiments to test against the influence that *each class* of complexity has on intelligence separately.

3.7 Implementation Details and Experimental Protocol

In this section I discuss some important test functionalities and experimental parameters. I also give a technical description of some example agent behaviours by showing how they can be practically evaluated over the Λ^* environment.

3.7.1 Setup and test parameters

The intelligence test was implemented in the C++ Object Oriented Programming language. The source code and scripts to run experiments have been released as open-source in (Chmait, 2016), with good efforts made to facilitate their re-usability.

Once the test is compiled and run, a new experiment is initiated. The number of test episodes ω , as well as the number of iterations in each episode, for that experiment can be entered into the command-line. Setting ω to 1000 episodes (runs) usually records a very small standard deviation between the test scores⁸. The size of the environment (and thus

⁸Usually a standard deviation of less than 0.001 is recorded between identical experiments.

the search space uncertainty $\mathcal{H}(\mu)$ [Recall Section 3.6.3]) as well as the number of agents to be evaluated can also be selected prior to each experiment. The robustness of the test scores depends on the size of the environment so it might be desirable to select a larger value of ω for larger environment spaces.

In each episode, agents are administered over different pattern complexities $K(\mu)$ automatically generated by the test, such that $K(\mu) \in [2, 23]$, where a $K(\mu)$ of 23 corresponds to, more or less, complex pattern prediction or recognition problems. Moreover, in each episode, the evaluated agents are automatically re-initialised to different spatial positions in the environment. (This will be revisited in Section 3.10). At the end of each experiment the (intelligence) scores (in the range $[-1.0, 1.0]$) of the evaluated agents and collectives, averaged over all test episodes, are displayed on the screen and also saved to file.

Agents can be evaluated in isolation as well as collectively following the agent environment framework previously described in Section 3.5. For instance, the test provides us with three key methods implementing the (multi) agent-environment framework. Let $\tilde{\mu}$ be an instance of the test environment Λ^* and Π a set of agents to be evaluated. The methods $sendObservations(\Pi, k)$ and $sendReward(\Pi, i)$ could be invoked on $\tilde{\mu}$ at each iteration i of the test in order to send observations and rewards respectively to all agents in Π , where $k \in \mathbb{N}^+$ refers to the k^{th} -Moore neighbourhood used to determine the breadth of the evaluated agents' observation range. At each iteration of the test, after receiving an observation, each agent in Π invokes its own $performAction()$ (polymorphic) method which returns a discrete action in the range $[1, 9]$, such that an action in $\{1, 2, 3, 4, 5, 6, 7, 8, 9\}$ maps position-wise to $\{up-left, up, up-right, left, stay, right, down-left, down, down-right\}$. The selected action is subsequently used to update the agent's position in the environment.

Despite the test's being originally designed to return a general measure of intelligence, I do not make this assumption in this chapter. Nevertheless, one can appraise the Λ^* environment, at a minimum, as an accurate measure of the testee's ability of performing over a class of: inductive inference, compression⁹ and search problems, all of which are

⁹I give here references to the Minimum Message Length (MML) principle (Wallace and Boulton, 1968), which has some relevance to this work and to the notions of intelligence, compression and inductive inference in general. Preceding the Minimum Description Length (MDL) (Rissanen, 1978) principle by ten years, Minimum Message Length (MML) (Wallace and Boulton, 1968; Boulton and Wallace, 1969, 1970; Boulton, 1970; Boulton and Wallace, 1973c,b,a, 1975; Wallace and Boulton, 1975; Boulton, 1975) has many applications. Some of these applications include angular data (Wallace and Dowe, 1993, 1994a; Dowe et al., 1996), mixture modelling of multinomial and Gaussian variables (Wallace and Boulton, 1968; Wallace, 1986) and extending this to also include *Poisson* and *von Mises* circular distributions (Wallace and Dowe, 1994b, 2000) and also some other distributions (Agusta and Dowe, 2002, 2003a,b), single latent factor analysis (Edwards and Dowe, 1998) and models allowing for sequential modelling (Edgoose and Allison, 1999; Edgoose et al., 1998) and spatial correlation (Wallace, 1998; Visser and Dowe, 2007), and other mixture modelling work (Figueiredo and Jain, 2002; Kasarapu and Allison, 2015)—as well as Bayesian nets with both continuous and discrete attributes (Comley and Dowe, 2003, 2005) (and see also Sec. 6.7 of this thesis). Also worth mentioning is work on MML decision trees (Wallace and Patrick, 1993), decision graphs (Oliver and Wallace, 1991; Oliver, 1993; Oliver et al., 1992; Tan and Dowe, 2002, 2003) and oblique decision trees (Tan and Dowe, 2004), and (e.g.) work on time series (Fitzgibbon et al., 2004; Schmidt, 2008), econometric panel data and estimation (Dowe, 2011, Sec. 6.5), learning generative models for structural representations (Torsello and Dowe, 2008a,b) and MML hypothesis testing (Dowe, 2008a, Sec. 0.2.5, p. 539 and Sec. 0.2.2, p. 528, col. 1 and Sec. 1)(Dowe, 2008b, p. 433 (Abstract), p. 435, p. 445 and pp. 455-456)(Musgrave and Dowe, 2010)(Dowe, 2011, Sec. 3.2, p. 919 and Sec. 7.6, p. 964)(Makalic and Schmidt, 2011). Comparisons of MML and MDL include, e.g., (Viswanathan et al., 1999; Fitzgibbon et al., 2004). For further references on MML—including comparisons with MDL—see, e.g., (Wallace, 2005; Dowe, 2011) and references therein. For information-theoretic work related to quantifying creativity, see

particularly related to intelligence (Dowe and Hajek, 1997a,b, 1998; Hernández-Orallo and Minaya-Collado, 1998; Sanghi and Dowe, 2003; Dowe et al., 2011). Note, however, that I will use the term *intelligence* to describe the effectiveness, or the accuracy, of an evaluated agent over this test. It is of great importance that the illustrative class of problems assessed by the test is shared across, and applies to, various types of cognitive systems since this meets my criteria for the evaluation, as raised at the beginning of this chapter (refer to Sections 3.2 and 3.3).

3.7.2 Modularity and code re-use

I have provided a large set of functionalities which might come in handy when amending and extending the current scope of the test, and for defining new agent behaviours to be evaluated. These functionalities can be found in (Chmait, 2016) as part of the utility class *General*, under the directory `/src/General.cpp`. Moreover, I have used `UnitTest++` (Llopis and Nicholson, last accessed, April 2016), a lightweight unit testing framework for C++ over Windows, in order to allow for easy defect isolation, assist in validating existing and newly implemented functionality, and encourage code review.

3.7.3 Implementing new agent behaviours

I defined an abstract class *Agent* that includes many useful functionalities. This makes implementing and evaluating new agent behaviours over the Λ^* environment faster, easier and less error-prone.

New isolated (non-interactive) agent behaviours can be introduced as (one of the) subclasses of *Agent*, providing implementations for its abstract methods as necessary. Interactive agent behaviours, on the other hand, are polymorphic classes redefining and extending the behaviour of their isolated agent's superclass.

Homogeneous collectives of interactive agents are aggregations of two or more interactive agents of the same behaviour (class). A simplified Unified Modeling Language (UML) class diagram illustrating the hierarchical relationships between isolated and collective agent behaviours is given in Figure 3.5. Likewise, heterogeneous collectives of agents can be defined as aggregations of two or more interactive agents of different behaviours (classes). Definitions of isolated and collective agent behaviours used for my experiments are given in the following section.

3.8 Agent Types and Behaviours

In this section, I give a (formal) description of the different agent behaviours that will be evaluated over the Λ^* environment following (Chmait, Dowe, Li, Green and Insa-Cabrera, 2016).

(Dowe, 2013, Sec. 4.6); and for other related information-theoretic work originally inspired by intelligence testing, see (Dowe, 2008a, footnote 175)(Dowe, 2008b, pp. 437–438)(Dowe et al., 2011, Sec. 3)(Dowe, 2013, Sec. 4.1).

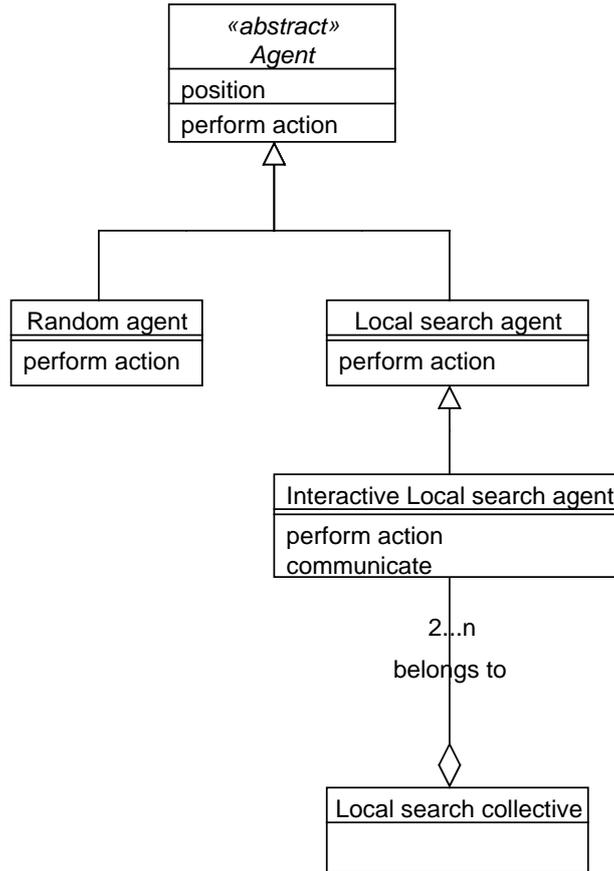


Figure 3.5: A simplified UML class diagram illustrating the relationships between (some) classes of isolated agent behaviours and agent collectives. A collective is an aggregation of two or more (objects belonging to one or more classes of) isolated agents (Chmait, Li, Dowe and Green, 2016, Fig. 4).

3.8.1 Local search agent

Given an agent π_j , I denote by c_j^i and $r(c_j^i)$ the cell where π_j is located at iteration i , and the reward in this cell respectively. Let N_j^i and $R(N_j^i)$ denote respectively the set of (Moore) neighbour cells of agent π_j (including c_j^i) at iteration i , and the reward values in these cells. $R(c_j^i, a)$ is a function that returns the reward agent π_j gets after performing an action $a \in \mathcal{A}$ when it is in cell c_j^i . The behaviour of a local search (\mathcal{LS}) agent π_j at iteration i is defined as follows:

$$a_j^i \leftarrow \arg \max_{a \in \mathcal{A}} R(c_j^i, a).$$

If all actions return an equivalent reward, then a random action in \mathcal{A} is selected.

3.8.2 Reinforcement learning agents

Two of the most frequently used Reinforcement Learning (\mathcal{RL}) behaviours are Q-learning (Watkins and Dayan, 1992) and Sarsa (Rummery and Niranjan, 1994; Watkins and Dayan,

1992). In the Q-learning behaviour, agents learn using an action-quality function in order to find the best action-selection policy for a given MDP (Markov Decision Process). Alternatively, Sarsa agents learn a MDP policy using an on-policy temporal-difference learning technique.

Q-learning agent

In this reinforcement learning behaviour, the evaluated Q-learning (Watkins and Dayan, 1992) agent learns using a state-action pair quality function, $Q : S \times \mathcal{A} \rightarrow \mathbb{R}$, in order to find the action-selection policy that maximises its rewards.

Each test episode of ϑ iterations is equivalent to one training session. Because the testing environment is dynamic, I define a Q-learning state $s_i \in S$, that an agent π_j occupies at iteration i , as the pair $\{c_j^i, i\}$ consisting of π_j 's current cell position c_j^i and the current iteration i , thus leading to a total number of states $|S| = m \times n \times \vartheta$, in a m-by-n environment space, over one test episode.

Before learning starts, the elements of the Q-table are initialised to 2.0 so that the quality of a state-action pair, $Q \leftarrow S \times \mathcal{A}$, remains positive despite that rewards fall in the range $[-1.0, 1.0]$. This might also promote the exploration of more state-action pairs before convergence. The \mathcal{RL} agents are trained for 100 rounds previous to each episode of evaluation using both a discount factor γ and a learning rate α of 0.30, selected after fine-tuning these parameters¹⁰ on a single agent scenario to reach a general (average) optimal payoff. The agents learn offline. Thus they cease to update their Q-table once their training is complete.

The Q-Learning behaviour over one training session is illustrated in Algorithm 1 using the notations from Section 3.8.1. During evaluation, the evaluated Q-learning agent travels

Algorithm 1 Q-Learning agent behaviour over one training session (Chmait, Li, Dowe and Green, 2016, Algorithm 1).

```

1: Initialize: Q-table, learning rate  $\alpha$  and discount factor  $\gamma$ .
2: Begin
3:   for iteration  $i \leftarrow 0$  to  $\vartheta - 1$  do                                ▷ loop over iterations
4:      $s_i \leftarrow \{c_j^i, i\}$                                              ▷ set current state
5:     execute  $a_j^i \leftarrow \arg \max_{a \in \mathcal{A}} Q(s_i, a)$                     ▷ perform action
6:      $s_{i+1} \leftarrow \{c_j^{i+1}, i + 1\}$                                 ▷ set new (post-action) state
7:      $Q(s_i, a_j^i) = Q(s_i, a_j^i) + \alpha \left[ R(c_j^i, a_j^i) + \gamma \max_{a \in \mathcal{A}} Q(s_{i+1}, a) - Q(s_i, a_j^i) \right]$   ▷ update
       Q-table
8:   end for
9: End

```

between states by performing the actions (leading to the states) returning the highest reward values recorded in its Q-table at the end of its training session.

¹⁰The discount factor $\gamma \in [0, 1]$ determines the importance of future rewards such that a γ value close to zero prompts the agent to consider immediate rewards (*myopic* agent), while a γ value closer to one assigns greater weight on future rewards and in that case the agent will strive for long-term high rewards. The learning rate $\alpha \in [0, 1]$ controls the speed of convergence of the Q-table by determining to what extent the newly acquired information about a given (environment) state will override the old information.

Sarsa agent

SARSA (short for State-Action-Reward-State-Action) (Rummery and Niranjan, 1994; Watkins and Dayan, 1992) is an alternative behaviour from Q-learning in which *the value of the policy the agent is actually carrying out* is learnt in such a way that it can be iteratively improved by taking into account the costs associated with exploration. In this thesis I refer to SARSA as Sarsa. The Sarsa behaviour over one training session is illustrated in Algorithm 2.

Algorithm 2 Sarsa agent behaviour over one training session.

```

1: Initialize: Q-table, learning rate  $\alpha$  and discount factor  $\gamma$ .
2: Begin
3:   for iteration  $i \leftarrow 0$  to  $\vartheta - 1$  do                                ▷ loop over iterations
4:      $s_i \leftarrow \{c_j^i, i\}$                                              ▷ set current state
5:     execute  $a_j^i \leftarrow \arg \max_{a \in \mathcal{A}} Q(s_i, a)$                  ▷ perform action
6:      $s_{i+1} \leftarrow \{c_j^{i+1}, i + 1\}$                                ▷ set (post-action) state
7:     simulate  $a_j^{i+1} \leftarrow \arg \max_{a \in \mathcal{A}} Q(s_{i+1}, a)$          ▷ simulate action
8:      $Q(s_i, a_j^i) = Q(s_i, a_j^i) + \alpha \left[ R(c_j^i, a_j^i) + \gamma Q(s_{i+1}, a_j^{i+1}) - Q(s_i, a_j^i) \right]$ 
       Q-table                                                               ▷ update
9:   end for
10: End

```

The same initialisation and training settings used for the Q-learning agents are also used for Sarsa agents. As mentioned previously, Sarsa agents learn a MDP policy using an on-policy temporal-difference learning technique.

3.8.3 Expert (oracle) agent

An expert or oracle agent knows the future movements of the special object \oplus . At each step i of an episode this agent approaches the subsequent $i + 1$ cell destination of \oplus seeking maximum payoff. However, if \oplus has a constant movement pattern (e.g., moves constantly to the right) pushing it away from the oracle, then the oracle will move in the opposite direction in order to intercept \oplus in the upcoming test steps. Once it intercepts \oplus , it then continues operating using its normal behaviour.

3.8.4 Random agent

A random agent randomly chooses an action from the finite set of actions \mathcal{A} at each iteration until the end of an episode.

The scores of the random and oracle agents are important for my experiments. These scores could be used as a baseline for the intelligence test scores of artificial agents where a random agent is used as a lower bound on performance while the expert agent is used as an upper bound.

3.9 Communication Protocols

Each of the agent types described in the previous section was also evaluated collectively as one unit (or group). The details of the interaction and communication protocols used by these agent collectives are given below.

3.9.1 Stigmergy or indirect communication

In Chapter 2 I looked into many areas of study in which collective intelligence was investigated. Many of the surveyed works in these areas were linked to the behaviour of animals and insects—swarms in particular. One of the main characteristics of complex swarming behaviours was the use of a form of indirect communication known as *stigmergy* (Grassé, 1959). An example of this behaviour can be observed in (the indirect coordination of) ants or termites (recall Section 2.2.1 and Table 2.1).

I propose a simple algorithm for enabling communication between local search agents using stigmergy (indirect communication/coordination). For instance, I let the agents induce fake rewards in the environment, thus indirectly inform neighbour agents about the proximity of the special objects. Note that fake rewards will not affect the score (real reward payoff) of the agents. Let $\hat{R}(N_j^i)$ denote the set of fake rewards in the neighbour cells of agent π_j (including c_j^i) at iteration i , and $\hat{R}(c_j^i, a)$ is a function returning the fake reward agent π_j gets after performing action $a \in \mathcal{A}$ when it is in cell c_j^i at iteration i . Fake rewards are induced in the environment according to Algorithm 3. Each agent proceeds

Algorithm 3 Stigmergic or indirect communication: fake reward generation over one iteration i of the test (Chmait, Dowe, Li, Green and Insa-Cabrera, 2016, Algorithm 1).

```

1: Input:  $\Pi$  (set of evaluated agents),  $0 < \gamma < 1$  (fake reward discounting factor), a test
   iteration  $i$ .
2: Initialize:  $\forall \pi_j \in \Pi: \hat{R}(N_j^i) \leftarrow 0.0$ .
3: Begin
4:   for  $j \leftarrow 1$  to  $|\Pi|$  do ▷ loop over agents
5:      $r^{max} \leftarrow \max R(N_j^i)$ 
6:      $r^{min} \leftarrow \min R(N_j^i)$ 
7:      $\hat{r} \leftarrow \gamma(r^{max} + r^{min})$  ▷ average expected reward
8:      $\hat{R}(N_j^i) \leftarrow R(N_j^i) + \hat{r}$ 
9:   end for
10: End

```

by selecting an action by relying on the fake rewards this time, instead of the real rewards, as follows: $a_j^i \leftarrow \arg \max_{a \in \mathcal{A}} \hat{R}(c_j^i, a)$. If all actions are equally rewarding, then a random action is selected. Thus, it is expected that local search agents using stigmergy to form non-strategic coalitions after a few iterations of the test as a result of tracing the most elevated fake rewards in the environment.

3.9.2 Implicit leadership through auctions and bidding

The processes of auctions and bidding are frequently observed in several contexts of our everyday lives from the simple sales of collectibles, antiques and art items, to the more complex commodity auctions at the level of corporations. Many types of auctions exist.

In this section, I define a simple cooperative setting where local search agents go into a *single dimensional English auction* (Parsons et al., 2011). At each iteration i , the agents bid openly against each other on the right to lead the other agents in their group by appointing one target cell to be approached. Moreover, at each iteration, each auctioneer (agent) generates a value of the maximum reward existing in its neighbourhood, which is then used as its bidding “money” for the auction. The richest agent¹¹ wins the auction visibly to all the other agents. The winner agent then selects the *target cell* to be approached by all other agents in the collective. This bidding behaviour is described in Algorithm 4 in which $n_j^i \in N_j^i$ and $r(n_j^i)$ denote one of the *Moore* neighbour cells of agent π_j (without excluding c_j^i) at iteration i and the reward in this cell respectively.

Algorithm 4 Single dimensional English auction at one iteration i of the test (Chmait, Dowe, Li, Green and Insa-Cabrera, 2016, Algorithm 2).

```

1: Input:  $\Pi$  (set of evaluated agents),  $-1.0 < \text{bid} < 1.0$ , a test iteration  $i$ .
2: Initialize:  $\text{bid} \leftarrow -1.0$ 
3: Begin
4:   for  $j \leftarrow 1$  to  $|\Pi|$  do ▷ loop over agents
5:      $\text{money} \leftarrow \max R(N_j^i)$ 
6:     if  $\text{money} \geq \text{bid}$  then
7:        $\text{bid} \leftarrow \text{money}$ 
8:        $\text{target} \leftarrow \arg \max_{n_j^i \in N_j^i} r(n_j^i)$  ▷ set the target to the neighbour cell  $n_j^i$  holding
           the highest reward  $r(n_j^i)$  at iteration  $i$ 
9:     end if
10:  end for
11: End

```

3.9.3 Imitating super-solver agents

Imitation is a phenomenon long observed in both human and non-human animals. For example mammals copy the actions of older members of their species, and so do human babies who imitate (the activities and impressions of) their parents. Imitation has also been used as a technique to help develop/design the behaviour of some artificial systems. All three types of cognitive systems (human, animal and machine) can learn by imitation and social learning (Fiorito and Scotto, 1992; Nehaniv and Dautenhahn, 2007).

I evaluated a group of *isolated local search agents* that is put in the same space with one (unevaluated) expert or oracle agent. Local search agents imitate the action of the oracle by following it into the same cell only when the oracle is within their visibility range

¹¹If more than one agent are equally rich then, for the sake of simplicity, the last one to participate in the auction wins.

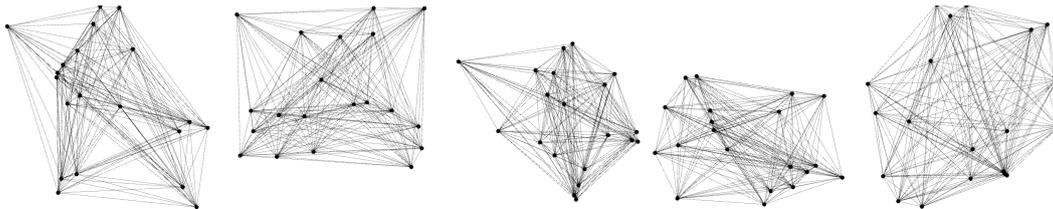


Figure 3.6: Examples of initial group topologies for collectives of 20 agents similar to the ones used in our experiments. The nodes represent the agents’ (cell) positions in the testing environment. Edges are added to improve the visibility of the topologies and simplify the comparison between them. The initial positions of the agents in the testing environment are generated using a *random* permutation (sampling with replacement) to ensure that the agents’ scores are valid regardless of their initial setup (group topology) in the environment space.

(consisting of their range of 1 Moore neighbourhood). Otherwise, the local search agents operate using their normal behaviour that was described earlier in Section 3.8.1.

3.9.4 Harnessing the wisdom of the crowd by information aggregation

Perhaps one of the most discussed phenomena perceived in groups and societies is the *wisdom of the crowd* (Surowiecki, 2005). Examples of this phenomenon were given in Section 2.2.2.

I implemented two cooperation mechanisms harnessing the wisdom of the crowd for both reinforcement learning and local search agents. In this cooperative scenario the opinion of all the evaluated agents is aggregated to form the collective opinion of the group.

In the case of reinforcement learning collectives, I let their members share and update a common Q-table, thus making them all learn and coordinate simultaneously. I evaluated both Q-learning and Sarsa collectives independently.

In the case of local search collectives, the observations of all agents in the collective are aggregated into one global observation (and rewards from these observations are averaged in the case of overlap). Then, each member proceeds by selecting the action maximising its reward in line with the global observation.

3.10 Experimental Setup

Before I discuss the results from the intelligence tests, I give a detailed description of the experimental setup I used to run my tests. Each experiment consists of 1000 episodes (runs) of the test, each episode in turn consisting of a number of iterations equal to 50. In each episode, agents are administered over a different task complexity $K(\mu)$, such that $K(\mu) \in [2, 23]$, where a $K(\mu)$ of 23 corresponds to a, more or less, complex pattern prediction or recognition task. Moreover, in each episode, the collectives are re-initialised with different topological (network) arrangements between their members. Examples of these topologies are illustrated in Figure 3.6.

Test experiments were conducted over different search space uncertainties $\mathcal{H}(\mu)$, and the (intelligence) scores (in the range $[-1.0, 1.0]$) of the evaluated agents/collectives averaged over the 1000 episodes were recorded. The score of the collective is calculated as the mean of the scores of its members. Recall that the metric of (individual agent) universal intelligence defined in (Hernández-Orallo and Dowe, 2010, Definition 10) was extended into a collective intelligence metric (Definition 3) returning an average reward accumulation per-agent measure of success (Definition 2) for a group of agents Π , over a selection of Λ^* environments.

I repeat here for convenience that local search agents were evaluated in isolation as well as collectively using four communication or interaction protocols:

- indirect coordination using stigmergy (Section 3.9.1),
- implicit leadership through auctions and bidding (Section 3.9.2),
- imitation (of the oracle/expert agent) (Section 3.9.3), and
- harnessing the wisdom of the crowd (*WOC*) through information aggregation (Section 3.9.4).

Likewise, reinforcement learning agents were evaluated in isolation and collectively by harnessing the wisdom of the crowd through sharing and updating a common Q-table (Section 3.9.4, last 2 paragraphs).

3.11 Results and Discussion

Table 3.1 shows sample results from the above-described experiments that are run over different environment (search space) uncertainties $\mathcal{H}(\mu)$. Results are given for both (isolated) individual agents and (interactive) collectives. The number of agents (or members) in a group used is $|\Pi| = 10$ agents. The standard deviation σ of the test scores illustrated in Table 3.1 is less than 0.001 between identical experiments. The results from Table 3.1 are plotted in Figure 3.7 to allow for analysis.

3.11.1 Collectives outperform individuals

Results from Figure 3.7 clearly show (in at least three separate cases) that cooperative or interactive individuals can be more effective than isolated ones. (From Definition 2, the score of the whole is more than the sum of its parts.) This is consistent with earlier results (e.g., (Mason and Watts, 2012; Bettencourt, 2009)) for obvious reasons owing to diffusion of information (e.g., synergy, information sharing, etc.) leading to the reduction of uncertainty inside the collective.

Yet the question remains, what are the dynamics that influence or control such results? By revisiting the main goals and objectives of this thesis (outlined in Section 1.3) we see that, besides the measurement of intelligence (Objective *Obj01*), the main aim from (**G01**) consists of quantifying and analysing the influence of a list of factors on (individual and collective) intelligence (Objectives *Obj02* and *Obj03*). I address each of these factors in detail in the remainder of this chapter.

	$\mathcal{H}(\mu)$ value in bits	13.2	15.6	17.2	18.5	19.6
1	Random agent (Section 3.8.4)	-0.00079	0.00048	0.00008	-0.00013	0.00002
2	Local search (\mathcal{LS}) agent (Section 3.8.1)	0.3365	0.1696	0.0936	0.0575	0.0423
3	\mathcal{LS} collective using stigmergy (Section 3.9.1)	0.4025	0.2555	0.1431	0.0829	0.0579
4	\mathcal{LS} collective harnessing the WOC (Section 3.9.4)	0.3828	0.3475	0.3118	0.2601	0.2110
5	\mathcal{LS} collective using implicit leadership (Section 3.9.2)	0.3744	0.2842	0.2143	0.1722	0.1438
6	\mathcal{LS} collective using imitation (Section 3.9.3)	0.5729	0.2880	0.1666	0.1022	0.0731
7	Q-learning agent (Section 3.8.2)	0.2516	0.0950	0.0484	0.0301	0.0207
8	Q-learning collective harnessing the WOC (Section 3.9.4)	0.4030	0.1832	0.0870	0.0482	0.0309
9	Sarsa agent (Section 3.8.2)	0.2708	0.1007	0.0501	0.0308	0.0228
10	Sarsa collective harnessing the WOC (Section 3.9.4)	0.4511	0.2042	0.1010	0.0563	0.0348
11	Oracle agent (Section 3.8.3)	0.8207	0.7905	0.7619	0.7339	0.7059

Table 3.1: Intelligence test scores for collectives of 10 agents across different environment uncertainties $\mathcal{H}(\mu) \in [13.2, 19.6]$ bits, evaluated for 50 test-iterations over the Λ^* environment (Chmait, Dowe, Li, Green and Insa-Cabrera, 2016, Tab. 1). A plot of these results is also found in Figure 3.7.

3.11.2 Communication and interaction protocol

We observe in Figure 3.7 that the effectiveness of the (same selection of) agents depends on the collective decision-making technique, or the communication protocol, used to aggregate the information received from these agents. For instance, adopting auctions in local search collectives to claim leadership can be more effective than using stigmergy over some settings.

Figure 3.7 also shows that, under certain circumstances, introducing heterogeneity in a group of local search agents by imitating a (super-solver) oracle agent leads to more effective coalitions that outperform their homogeneous (and isolated) peers by aggregating new information into the collective. However, the comparison between local search collectives is rather more complicated as their intelligence measures seem to further depend on the uncertainty of the testing environment, and not only on the interaction protocol. We also

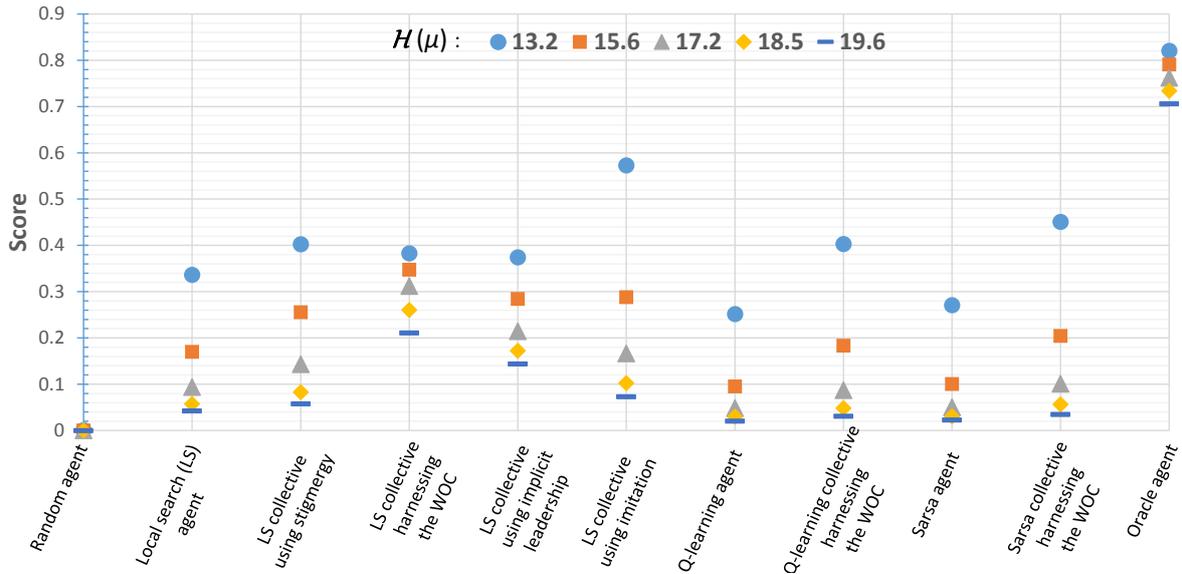


Figure 3.7: A plot of the test scores appearing in Table 3.1 (Chmait, Dowe, Li, Green and Insa-Cabrera, 2016, Fig. 3).

observe that harnessing the wisdom of the crowd by aggregating the observations of local search agents is very effective over highly uncertain environments, yet not exceptionally efficient in the opposite situation. The latter protocol seems to be very robust (in comparison to others) with respect to the changes in the uncertainty of the search space. A more thorough analysis on the efficiency of the examined communication protocols over different problem uncertainties is addressed in the following subsection.

In the case of reinforcement learning (\mathcal{RL}) agents, we observe that agents of different types (namely Q-learning and Sarsa) using the same cooperation technique to aggregate their information have achieved different scores. Sarsa agents outperform Q-learning agents, up to about a similar extent, in both cooperative and isolated settings. This indicates that, to a certain extent, the intelligence of the group also depends on (and is correlated with) the individual intelligence (or the type) of the agents in the group.

Overall, the scores in Figure 3.7 show that, despite the broad differences in the interaction protocols that are implemented and the wide range of task complexities used for assessment, there is evidence of collective intelligence in all the evaluated collectives. This shows that collective intelligence can emerge in a non-human context or environment, thus reinforcing and adding to the conclusions drawn in (Engel et al., 2015).

3.11.3 Uncertainty in the environment

Figure 3.7 shows that the performance of the evaluated agents decreases with the increase in uncertainty¹² $\mathcal{H}(\mu)$, in accordance with former tests (Insa-Cabrera, Dowe, España-Cubillo, Hernández-Lloreda and Hernández-Orallo, 2011) that have been applied on humans and artificial agents. Moreover, the gap between the scores of the isolated and

¹²Except for random agents, which always score around zero.

cooperative agents varies in view of the uncertainty in the environment, but the relationship between both variables cannot be easily grasped from the figure.

I wish to measure the variation in the weight of the cooperative agents' scores to their score in the isolated setting, across different environment uncertainties. Therefore, I define the *coefficient of effectiveness* $\theta = \alpha/\beta$, as the ratio of the score of a set of agents Π working in some cooperative scenario ($\alpha = \text{score}(\Pi^{\text{coop}})$), to its score in the isolated scenario ($\beta = \text{score}(\Pi^{\text{isolated}})$). I calculated θ for the different agent types, across different environment uncertainties $\mathcal{H}(\mu)$, and plotted the results in Figure 3.8.

Figure 3.8a shows that the *coefficient of effectiveness* θ values corresponding to local search groups using imitation ($\theta^{\text{imitation}}$) are more or less steady across the different $\mathcal{H}(\mu)$ values. A uniform $\theta^{\text{imitation}}$ value implies that local search agents using imitation are approximately equally more advantageous than (the same number of) isolated local search agents across the different problem uncertainties. This is also the case for local search agents relying on stigmergy ($\theta^{\text{stigmergy}}$). In addition, we observe that $\theta^{\text{imitation}} > \theta^{\text{stigmergy}}$ over the selected uncertainties, and thus collectives relying on imitation are more effective than those using stigmergy.

The observations are more interesting for collectives using auctions to claim leadership. For instance, in environments of uncertainties lower than 16 bits, imitating a smart agent is more advantageous than following a leader (bidding). Whereas, the inverse is true for environments of higher uncertainties. The effectiveness of local search agents using auctions significantly increases to become much higher than that of the same group of agents imitating an oracle. Similar results are observed for local search collectives harnessing the wisdom of the crowd. We conclude that relying on the best (super-solver) agent in the group does not guarantee an optimal performance. This is somewhat consistent with (Hong and Page, 2004)'s claims regarding diversity vs. ability—even though the intuitions here are different. These results have a fundamental impact on the choice of the communication protocol to be used in order to aggregate the information received from a group of agents, especially over problems where the search complexity can be estimated in advance.

For reinforcement learning agents, Figure 3.8b shows an overall similar shift in effectiveness for both Q-learning and Sarsa collectives. Their performances significantly increase over isolated agents to reach a peak around $\mathcal{H}(\mu) = 16$ bits, but then start to drop down over higher uncertainties. This illustrates the fact that cooperative reinforcement learning collectives are particularly advantageous over environments that are somewhat highly-uncertain for their isolated peers to perform efficiently in, yet not too uncertain for them (the cooperative collectives) such as to hinder their performance. In other words, collective intelligence is hardly perceived in (i) groups operating in very simple environments where individuals could perform relatively well, or in (ii) those environments that are too difficult (broad) to be explored within a limited *interaction time*, given a limited *number of members* in the group.

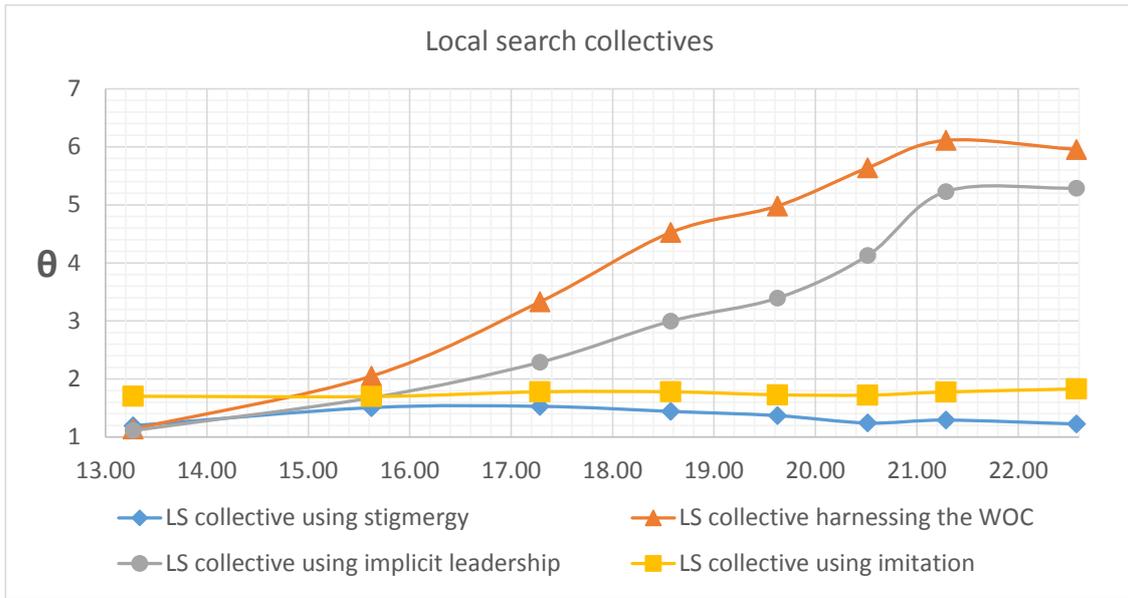
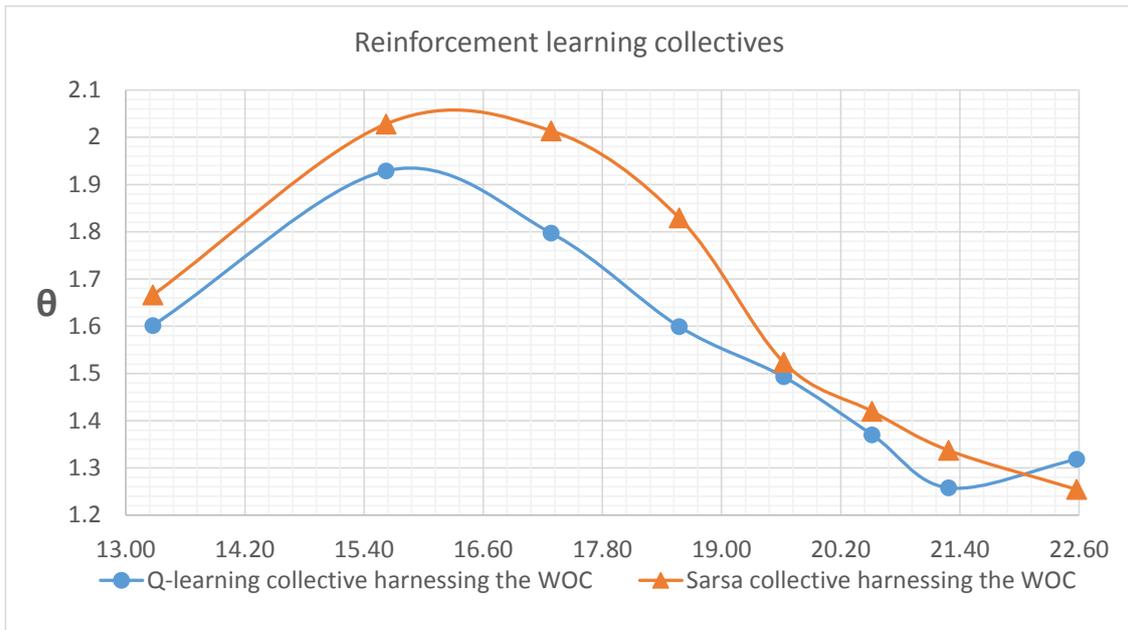
(a) Local search (\mathcal{LS}) collectives.(b) Reinforcement learning (\mathcal{RL}) collectives.

Figure 3.8: Shift in effectiveness θ for local search and reinforcement learning (\mathcal{RL}) agents over different environment uncertainties measured in bits (Chmait, Dowe, Li, Green and Insa-Cabrera, 2016, Fig. 4). Note that the standard deviation of the test scores is less than 0.001 between identical experiments, with each experiment consisting of 1000 runs (episodes) as described in our experimental setup in Sec. 3.10.

In order to understand the global picture of the agents' collective behaviour and its dynamics, it is crucial to look into the latter two factors (interaction time and number of members in a group) and measure their effects, if any exist, on group intelligence.

3.11.4 Number of agents in a group

In all my previous experiments the number of evaluated agents in a collective was set to 10 agents whereas Figure 3.9 illustrates the scores of the evaluated collectives across different number of agents varying between 5 and 70.

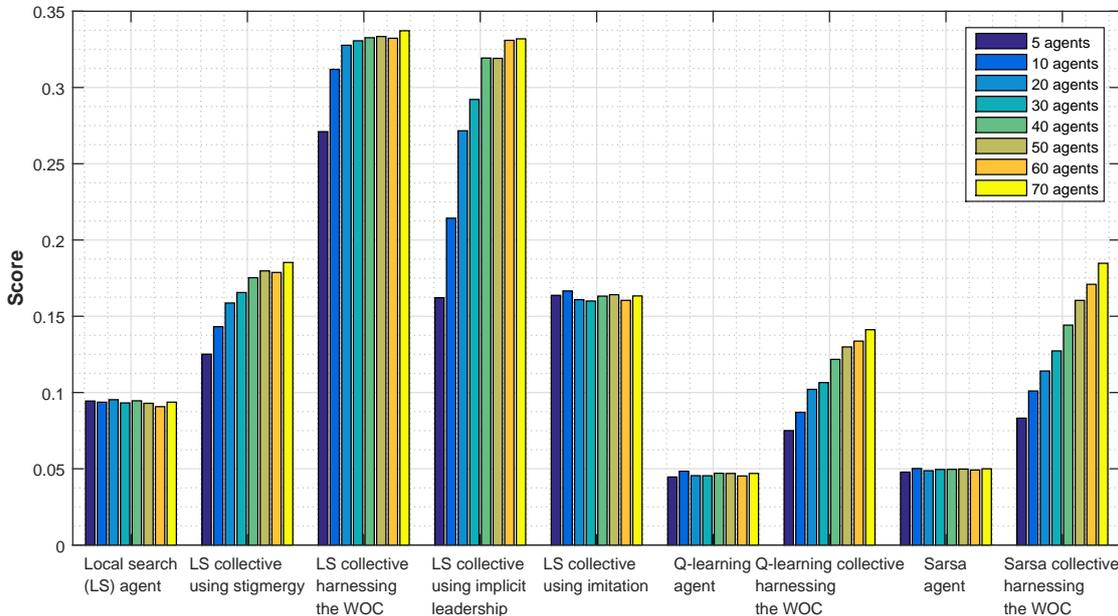


Figure 3.9: Intelligence scores recorded for collectives of different numbers of agents $5 \leq |\Pi| \leq 70$, in 17.8-bit $\mathcal{H}(\mu)$ environments. The scores of the isolated agents are also illustrated for comparison purposes (Chmait, Dowe, Li, Green and Insa-Cabrera, 2016, Fig. 5).

The general picture shows that local search collectives relying on stigmergy and auctions gradually improve in performance as more agents are added into the collective. This is not the case for collectives relying on imitation which only show a shallow variation in score. In fact, local search agents relying on imitation performed better than those using auctions when $|\Pi|$ was set to 5 agents. However, the opposite was true when we increased the number of agents to 10 and higher. This illustrates that, when the group is small in number, relying on a super-solver agent might be more advantageous than interacting between the individual members, however, as the group gets larger, more information is added into the collective and the expertise of a single oracle becomes rudimentary in comparison to the aggregated experiences (synergy) from individual members. Moreover, we observe that local search collectives harnessing the *WOC* improve faster in performance than those following a leadership. Nonetheless, when the number of agents gets higher the performances of these two collectives get closer to one another.

On one hand we observe that, increasing the number of local search agents is more effective, and has greater influence on the scores, for agents relying on auctions than those using stigmergy. On the other hand, the increase in efficiency is slightly non-linear to the number of agents introduced. For instance, the main improvements in scores are more concentrated at the early introductions of agents. Afterwards the scores continue to rise, but less and less significantly.

Similar observations illustrate that the effectiveness of reinforcement learning collectives harnessing the wisdom of the crowd improves as we increment the number of agents. Moreover, Sarsa collectives seem to be slightly more efficient than Q-learning collectives as new agents are introduced into the group. The key issue in this experiment is that, collective intelligence cannot be considered independently of the number of members in the group. Instead, it is a function of (so far) at least three factors, each having a different influence that we have measured, and bearing distinctive properties of which we have identified some.

3.11.5 Time and intelligence

In this paragraph I address the relevance of time to intelligence, which is often neglected in the assessment of collective intelligence¹³.

Figure 3.10 shows the variations in the intelligence scores as we extend the interaction time (number of iterations or interaction steps) of the test. We observe that some scores incline to converge as more time is given to the members to perform on the test. Figure 3.10a shows that the advantage of cooperative local search agent groups over isolated agents is higher at the early stages of the test in the case of agents using auctions to claim leadership. Afterwards, the gap in performance slowly decreases with time until iteration number 600. On the contrary, the gap between the scores of local search agents imitating an oracle and their isolated peers grows as we let the test run, implying that local search agents relying on imitation require longer periods of time to reach their best performance.

I have already shown in Figure 3.7 that, over some uncertainties, local search agents relying on auctions outperform those imitating an oracle. This is consistent with the results in Figure 3.10a (up to iteration 300). However, this experiment also suggests that imitating a super-solver agent is highly rewarding over time, leading to better-scoring collectives than when using leadership through auctions. These results illustrate how diverse social organisations between the members of the collective determine its performance over time. For instance, a (dynamic) leadership scheme or organisation seems to be more rewarding than a simple flat hierarchy relying on stigmergic communication given a limited interaction time with the environment.

Moreover, the general picture shows that a local search collective harnessing the *WOC* is most advantageous over isolated agents (and other collectives) mainly before the 300th iteration, at which point its performance begins to converge slowly.

In the case of reinforcement learning agents (Figure 3.10b), both isolated agents and collectives improve in performance with time, keeping an overall steady relationship between the differences in their scores. This raises another concern regarding the intelligence of artificial agents. It is intriguing as to what ideally counts as more intelligent, a fast reactive agent with a humble performance, or a slow one with an exceptional performance over an extended period of time? Should one consider the *potential intelligence* (Hernández-Orallo and Dowe, 2013) of an agent instead?

¹³Neglected issues of time and space resolution are highlighted in (Dowe, 2013, Sec. 4.4, p. 23-24).

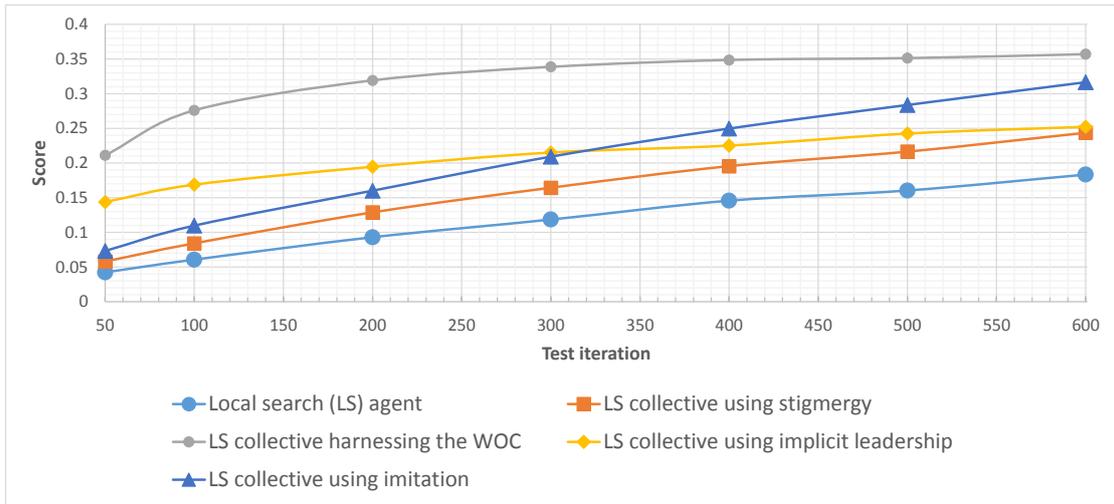
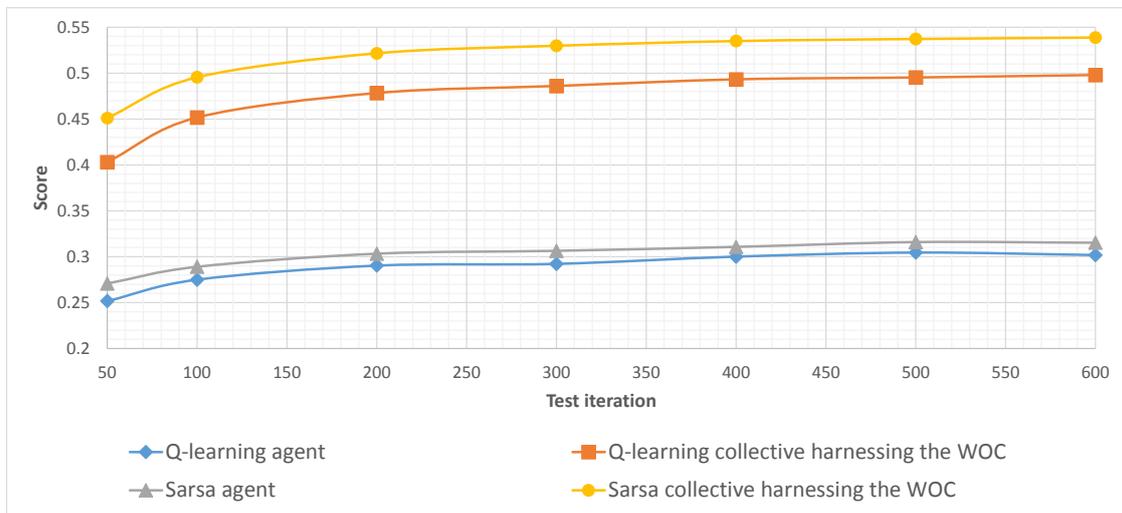
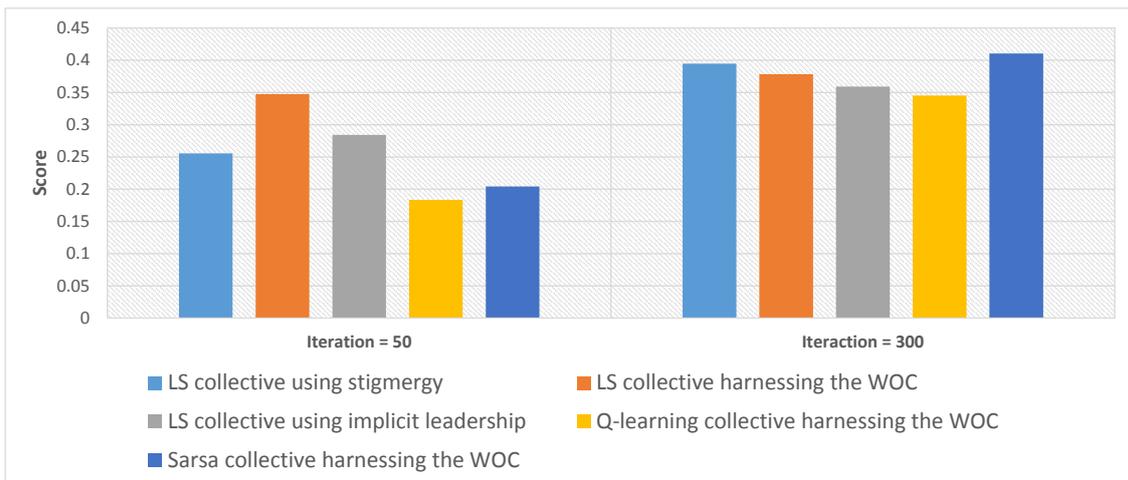
(a) Local search agent collectives in $\mathcal{H}(\mu) = 17.2$ bits.(b) \mathcal{RL} collectives in $\mathcal{H}(\mu) = 13.2$ bits.(c) Scores in $\mathcal{H}(\mu) = 13.2$ bits after 50 and 300 iterations.

Figure 3.10: Variations in intelligence scores as we extend the evaluation time (number of iterations) of the test (Chmait, Dowe, Li, Green and Insa-Cabrera, 2016, Fig. 6).

To understand the importance of time in measuring intelligence, I compare the scores of reinforcement learning and local search collectives over 13.2-bit $\mathcal{H}(\mu)$ environments after 50 and 300 test iterations as illustrated in Figure 3.10c. We find that local search collectives outscored Sarsa collectives up to the first 50 iterations while the opposite is true at iteration 300. This type of experiment is one of the most revealing of how the (communication and interaction) reasoning/learning speed of multiagent systems influences their measured performance given a finite/bounded operation or interaction time.

3.11.6 Algorithmic complexity and intelligence

In this paragraph I shed some light on how the performance of (groups of) agents is influenced by the algorithmic complexity of the task. To minimise the effect of search and exploration (relative to exploitation) on the scores, I initialised all agents to neighbouring locations from the \oplus special object. I then evaluated the agents over tasks of different algorithmic complexities (seeming randomness) $K(\mu)$ grouped into three difficulty levels: easy $\in [6, 8]$, medium $\in [9, 13]$ and hard $\in [14, 19]$. This experiment stands out from previous related experiments in the field, as collectives are assessed against tasks of quantifiable algorithmic complexities, as opposed to ones qualitatively ranked based of their difficulty.

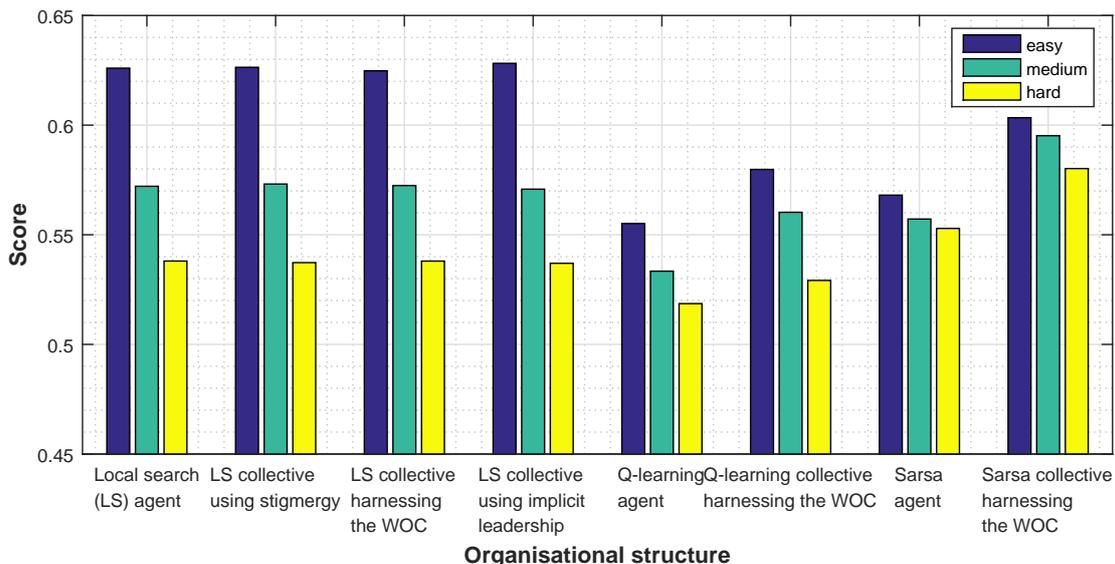


Figure 3.11: Scores over different task complexities $K(\mu)$ using collectives of $|\Pi| = 10$ agents, evaluated in 13.2-bit $\mathcal{H}(\mu)$ environments for 50 interactions (Chmait, Dowe, Li, Green and Insa-Cabrera, 2016, Fig. 7).

Results illustrated in Figure 3.11 show that the performance of artificial agents, similar to that of individual human performance (Insa-Cabrera, Dowe, España-Cubillo, Hernández-Lloreda and Hernández-Orallo, 2011), decreases when evaluated over patterns of higher algorithmic complexities. For instance, learning and predicting (seemingly) random patterns is more difficult, per se, than learning or inferring compressible ones.

Moreover, this experiment suggests that reinforcement learning collectives are better learners than their isolated peers since the difference between the cooperative agents' scores

over the three levels of difficulty is significantly smaller than that of the isolated ones. What's more intriguing in Figure 3.11 is the difference in behaviour between cooperative \mathcal{RL} agents and local search collectives. While \mathcal{RL} collectives are still more effective over isolated agents when evaluated over learning problems, all local search agents (isolated and collectives) performed equally when the effect of *search and exploration* was minimised. More importantly, we find that \mathcal{RL} collectives are more robust with respect to the change in algorithmic complexity as opposed to local search agents which display a wide gap in scores over the three levels of complexity.

All in all, what this experiment suggests is that, further to the previously examined factors, (collective) intelligence is a function of the agent type and the algorithmic complexity of the given task, both combined.

3.12 Organisational Behaviour

In spite of the different communication protocols I have evaluated, it is still not clear how the organisational structure of the group (Child, 1972; Mintzberg, 1979; Tran and Tian, 2013) affects its performance on intelligence tests. As discussed earlier in Section 2.2.2, this is an important factor that has increased the competitive advantage of many (online) businesses and shaped their management operations.

Therefore, I have further evaluated the performance of equally sized collectives of local search agents organised in four different (divisional and network) structures and studied their organisational behaviour. These structures are illustrated in Figure 3.12 and described below:

- In the flat, fully connected, structure (Figure 3.12a) all agents share their observations between one another. This absolute aggregation of information leads to a similar effect as that of local search collectives harnessing the wisdom of the crowd.
- In the subgroup structure (Figure 3.12b) I divide the collective into four smaller subgroups. Each one of those subgroups then implements a flat structure as the one described previously.
- In the hierarchical structure (Figure 3.12c), each (non-leaf) agent receives feedback from its children at each iteration of the test before performing an action. Leaf-nodes operate in isolation.
- Finally, in the autocratic structure (Figure 3.12d), a single agent controls the actions of the rest of the collective irrespective of its members' observations.

The results from my experiments show that flat, fully-connected, network structures are the most efficient since they maximise the aggregation of information received from the members of the collective. However, it is known that this type of structure is very costly as it requires a large number of connections¹⁴ to be introduced between the members of the

¹⁴The number of connections in a flat, fully-connected, network structure is equal to $n(n-1)/2$ connections, where n is the number of agents.

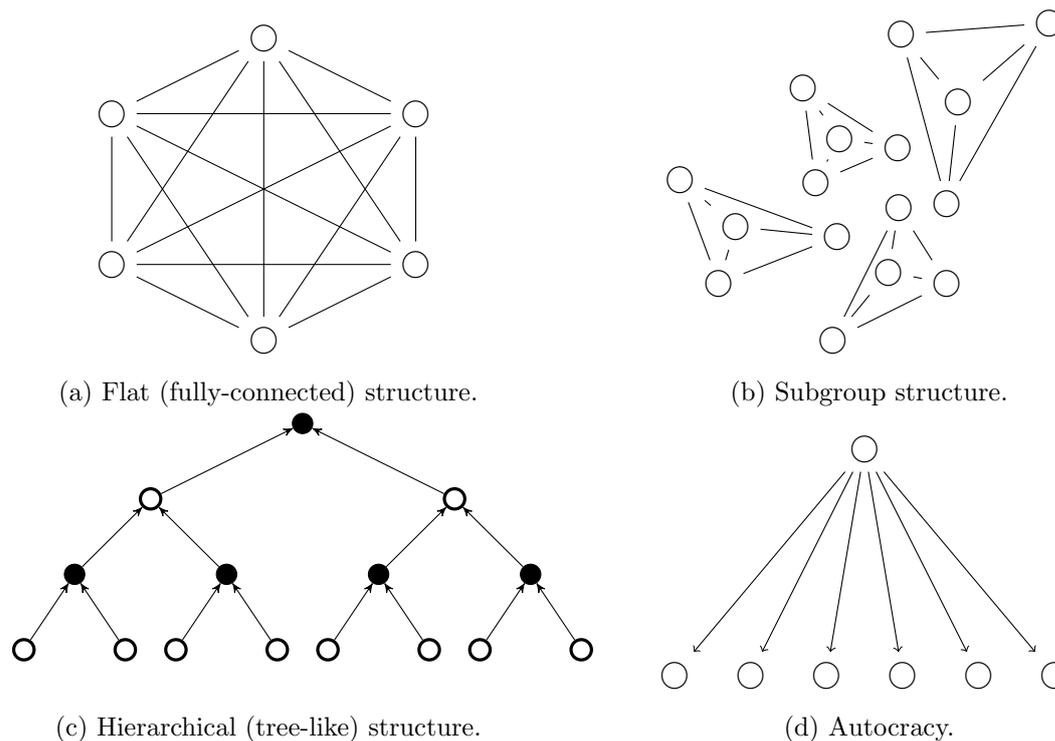


Figure 3.12: Graphical representation of different group organisational structures (or group topologies). Nodes represent agents and edges reflect the flow of communication and interaction between these agents. Undirected edges mean that communication can occur in both directions (Chmait, Dowe, Li, Green and Insa-Cabrera, 2016, Fig. 8).

group (Strogatz, 2001) which can be extremely difficult to implement and coordinate in real world organisations or businesses. In this type of organisation (and for the corresponding environment uncertainty and evaluation-time parameters) the number of agents did not significantly affect the performance of the collective. Whereas, after dividing this collective into smaller subgroups, the number of agents turns out to be of major importance¹⁵. The effectiveness of each subgroup improved gradually with the increase in the number of agents thus reducing the gap in performance between this organisational structure and fully-connected one. This shows that dividing a collective into smaller groups is most beneficial for highly populated collectives, especially when the number of connections inside the collective grows very large and becomes a bottleneck on communication.

In the hierarchical and autocratical structures, the measured effectiveness is low compared to the previous two models. We observe that the average performance of a hierarchical group is slightly steadier than that of a group governed by single agent with absolute control on decision-making. Interestingly, we also notice that in the hierarchical structure, high-scoring agents are the ones at the top of the hierarchy since (in this model of interaction) they receive feedback from their children (which in turn receive feedback from theirs) while the ones at the bottom perform in isolation and have low scores. Figure 3.13 shows that the average scores of the root agents in the hierarchy are significantly higher

¹⁵Note that all four evaluated subgroups showed a similar performance, but scores were only plotted for the first subgroup to enhance readability.

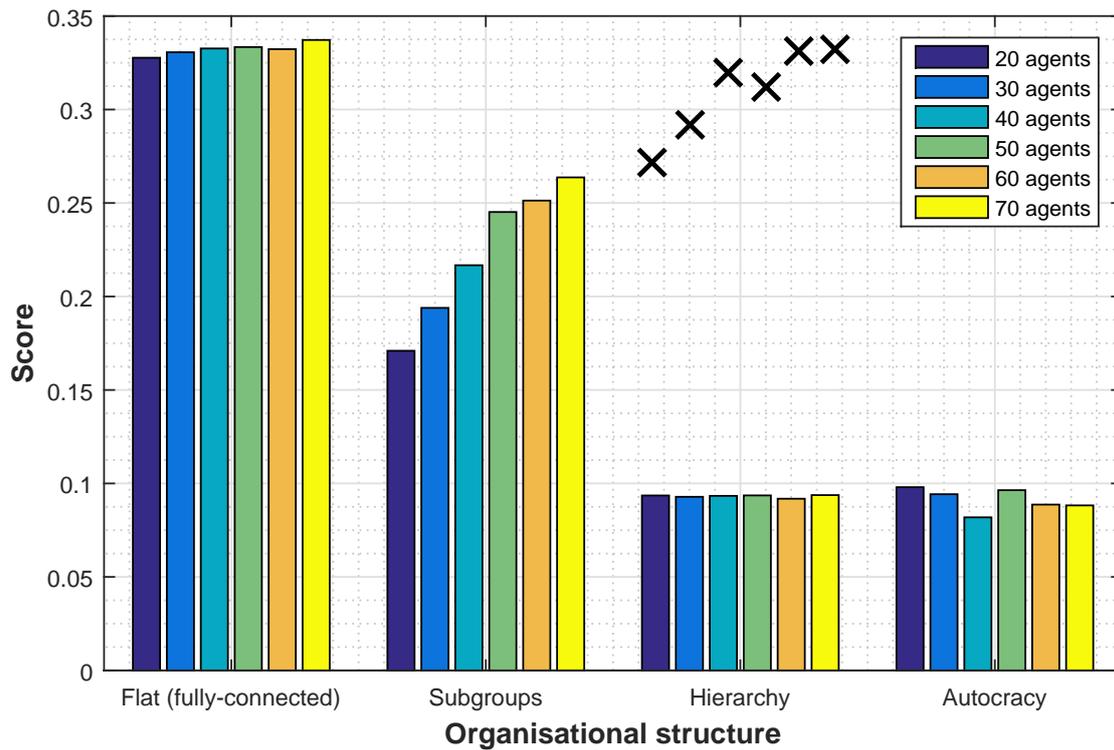


Figure 3.13: Scores of local search collectives organised in different network structures, across various number of agents between 20 and 70 (Chmait, Dowe, Li, Green and Insa-Cabrera, 2016, Fig. 9). The collectives are evaluated in 17.2-bit $\mathcal{H}(\mu)$ environments for 50 iterations. The label (X) in this figure depicts the average score of the agent at the top of the group hierarchy (that is the agent equivalent to the root node of the binary-tree used to represent that hierarchy).

than the average score of the collective, indicating a high standard deviation between the members' scores in this organisational structure. Since the number of leaves is almost half the number of nodes¹⁶, and the number of agents declines quickly as we move up the hierarchy, this organisation does not deliver a high average group performance.

Finally, Figure 3.13 shows that the performance of a local search collective implementing an autocracy is similar on average to that of isolated local search agents. Agents in this organisational structure do not show any significant discrepancies in their scores or behaviours.

3.13 Alternative Environments and Further Considerations

The Λ^* (Lambda Star) environment focuses on an important, but restricted, set of canonical tasks particularly relevant to the intelligence of AI agents. Nonetheless, the generalisation of these canonical tasks does not account for the complete range of multiagent problems. In particular, the tasks to perform in the Λ^* environment are a nice abstraction of two problems (among others) in the literature: searching for a moving target while avoiding injury, and nest selection when there is one and only one best nest. But these tasks do not capture other multiagent problems like those that require explicit coordination (e.g., lifting and moving a table).

3.13.1 Measuring multiagent coordination

Coordination is an important feature in multiagent systems which has a high influence on their performance. Measuring coordination between interactive agents can be a difficult task. The scope of the Λ^* (Lambda Star) environment does not currently account for the measurement of advanced forms of coordination between agents, but I am considering extensions to assess this. For instance, problems that require coordination can be evaluated if the payoff received from the *Good special object* \oplus (recall Section 3.6.1) had only occurred if two or more agents are in its neighbourhood at one iteration.

Another interesting extension to the test setting is to enable the environment to respond to the agent's behaviour and actions. Testing can be performed in an even more heterogeneous setting where the agents don't have the same reward function and/or actions and observations, and to give more attention or weight to the agent's learning (ability), which is an important aspect of intelligence.

Moreover, other properties like (environment) coverage could be evaluated by dispersing agents in the space to monitor what is happening in the environment (e.g., monitoring which neighbourhoods are/are not explored by the agents).

3.13.2 Fitness landscapes

The Lambda Star (Λ^*) is one of many environments which can be used to evaluate artificial agents. A famous problem in AI is to evaluate the performance of artificial agents over fitness landscapes consisting of many local optima but only one global optimum. The

¹⁶Number of leaves in a full binary tree is equal to $(\#nodes + 1)/2$.

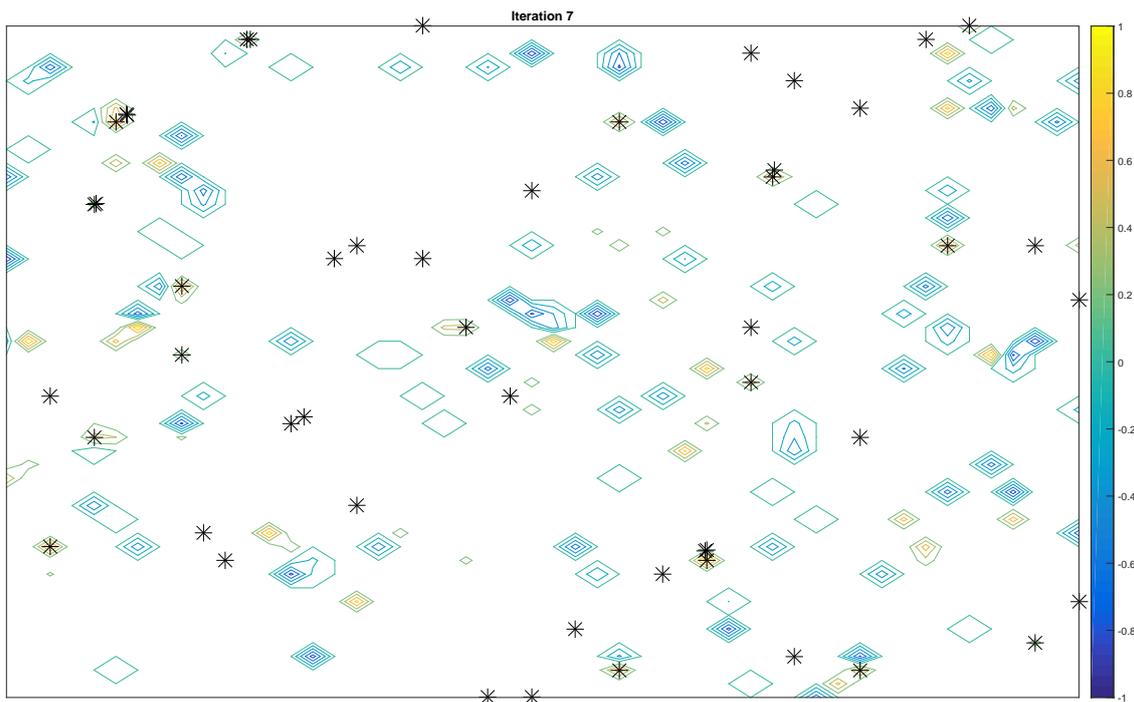


Figure 3.14: This is a screen-shot from the early stages of a (multiagent) simulation showing a fitness landscape with many local optima but only one global optimum. The colours (and their different intensities showing in the right-hand side colour-bar) represent fitness, or the quality of the landscape (ranging between $[-1.0, 1.0]$), at different coordinate/cell positions. The black stars represent the positions of the agents operating in the landscape (Chmait, Li, Dowe and Green, 2016, Fig. 5).

landscapes reflect the evaluations of some fitness (or utilisation) function over a set of candidate solutions. Adaptive landscapes can be considered where the underlying fitness evolves or changes over time.

I have implemented an extended version of the Λ^* environment to assess the performance of agents searching (dynamic) fitness landscapes where rewards and penalties are distributed over different areas of the landscape. I further designed a simulation depicting the behaviour of (co-operative) artificial agents exploring such landscapes over a period of time. This will be used for my future work in order to assess the trade-off between exploration and exploitation in a reinforcement learning setting, and investigate the influences of this trade-off on the agents' payoff in a multiple candidate solution space or environment. A screen-shot from the early stages of a sample simulation is given in Figure 3.14.

3.13.3 Further thoughts on robust intelligence tests

The same way collective intelligence can emerge between artificial agents (due, for example, to the wisdom of the crowd (Surowiecki, 2005), information sharing, reduction in entropy, etc.), pluralistic ignorance (Seeme and Green, 2016) is also a common phenomenon observed in many social settings which can occur between rational agents. Robust intelligence tests should be able to detect such a phenomenon. For instance, the field of game theory has highlighted several scenarios where cooperation between agents does not lead

to an optimal payoff (e.g., the famous prisoner’s dilemma (Poundstone, 2011)). A robust intelligence test should be general enough to reflect and evaluate such scenarios.

Other multiagent phenomena witnessed in various social settings reflect how agents acting individually might perform adversely to the common good, and thus deplete their available resources as a consequence of their collective behaviour. A robust intelligence test should allow for a quantitative assessment of *the tragedy of the commons* (Hardin, 1968) phenomena occurring in multiagent scenarios.

3.14 Conclusions and Future Work

I have addressed in this chapter the relevance of several factors and their interaction to the notion of intelligence and its emergence. I first started by looking at the different contexts in which collective intelligence has been shown to emerge, from face-to-face human groups, group collaborations via the web, social insect colonies and swarms, etc. Accordingly, I filtered a series of factors and features that are not coupled to one particular cognitive system, problem or environment, and illustrated how they influence the collective behaviour of the group, and hinder its intelligence.

The studied factors were shown to have a major influence on the performance of collectives—that I have also measured. But, I am most intrigued by the peculiar nature of collective intelligence seen as a function of all the examined factors simultaneously, as well as some of them combined. I identified circumstances where one cooperative system outperformed another under some values or setups of the studied factors yet failed to do so under others (e.g., in Section 3.11.5, limited vs. extended interaction time and, in Section 3.11.3, low vs. high environment uncertainty), reflecting on how these factors independently but also jointly shape the effectiveness of multiagent systems, and the spread of intelligence in these systems.

Some of my conclusions (in Section 3.11.3) reflected how relying on an expert (imitating a super-solver) agent in the group does not necessarily guarantee its optimal performance. I also measured the effect of introducing more agents into the group (Section 3.11.4), and showed that it is tightly controlled by the communication protocol used between its members. I have highlighted scenarios (in Section 3.11.6) where only some types of collectives outperform their equally sized group of isolated agents over (algorithmically) complex environments, and shown how the influence of the environment difficulty (uncertainty and complexity) is a major factor controlling the capacity for intelligence. Moreover, I looked (in Section 3.12) into how the effectiveness of (the same selection of) agents adopting different organisational and network structures can significantly vary from one structure to another.

This chapter tackled the first goal (**G01**) of this thesis by investigating some fundamental questions in AI and showing the existence, and *quantitatively* measuring the influence, of some general factors and principles shaping the spread of intelligence that are regularly perceived across different cognitive systems.

I have released the source code and scripts to run my experiments as open-source in (Chmait, 2016) to encourage additional testing and extensions to (the current version of)

the A^* environment. This will allow other researchers in the AI community to quantitatively evaluate new types of heuristics, algorithms, communication protocols and network structures. The motivation is to help transfer the results from this chapter into a guideline for understanding and designing multiagent cooperation.

Another future goal is to further evaluate agents and collectives over a wide range of general AI problems. For example, agents (isolated or collectives) could be evaluated over exploration/exploitation problems in an environment consisting of a hidden fitness landscape with many local, and only one global, optima. Other possible examples might include different forms of pattern recognition (and sequence completion) problems, in which payoff is determined by how accurately a subject learns and predicts a pattern. Other general multiagent problems that require coordination (e.g., lifting and moving an object), or scheduling (e.g., exam timetabling (Chmait and Challita, 2013)), can be used as alternative evaluation techniques.

Chapter 4

Observation Communication and Intelligence

Intelligence is the ability to adapt to change.

—Stephen Hawking, A Briefer History of Time (2005)

The research in this chapter has been published in the following article:

- Nader Chmait, David L. Dowe, David G. Green, and Yuan-Fang Li (2015). *Observation, communication and intelligence in agent-based systems*, in J. Bieger, B. Goertzel and A. Potapov (eds), Proceedings of the 8th International Conference on Artificial General Intelligence, Berlin, Germany, Vol. 9205 of Lecture Notes in Artificial Intelligence (LNAI), Springer, pp. 50-59.
http://dx.doi.org/10.1007/978-3-319-21365-1_6

4.1 Overview

The intelligence of multiagent systems is known to depend on the communication and observation abilities of its agents. Our experiments from the previous chapter revealed the existence of a significant influence of the factor of communication on the performance of agent groups. Nonetheless, not much has been said about the observation abilities of the agents. For instance, it is not clear which factor, observation or communication, has the greater influence on performance. By following a similar experimental approach to the one in Chapter 3, I will measure and compare the impact of each of these two factors, individually, on the intelligence of multiagent systems. In order to achieve this, I will present a method to—in turn—quantify the communication and observation abilities of an agent which, as we will see later on, I measure using the notion of Shannon’s entropy (Shannon, 1948). Some of the outcomes from this chapter indicate that the effectiveness of multiagent systems with low observation or perception abilities can be significantly improved by using high communication entropies between the cooperative agents in the system. I also identify circumstances where these assumptions fail, and analyse the dependency between the factors of observation and communication.

4.2 Introduction

The literature on multiagent systems has put forward many studies showing how factors such as *communication* (Dowe et al., 2011; Bettencourt, 2009; Panait and Luke, 2005) and *observation* (Fallenstein and Soares, 2014; Weyns et al., 2004; Franklin and Graesser, 1997) can influence the performance of multiagent systems. Our work from the previous chapter has also touched upon several topics inherent to the notion of intelligence testing in the context of cooperative multiagent systems. However, it is still ambiguous whether in a group of interactive agents:

- (i) augmenting the agents' observations to read and interpret the environment in which they operate, or rather

- (ii) boosting the communication between these agents,

has the higher influence on their performance. In fact, one of the fundamental characteristics of agent-based systems is their ability to observe/perceive and sense the environment (Wooldridge and Jennings, 1995; Franklin and Graesser, 1997). Within a multiagent system setting, perhaps the main property of agents is their ability to interact and communicate (Wooldridge and Jennings, 1995, Section 5).

In this chapter, I aim to analyse the factors of communication and observation by measuring and comparing the influence that each has on the intelligence of cooperative agent-based systems. Moreover, I will investigate and measure how these two factors are related to one another.

The motivations identified in (Section 3.3 of) Chapter 3 also apply here. Additional motivations particularly linked to the identified aims of this chapter also come to light. For instance, in real-world multiagent applications, agents can have limited sensitivity of their environment (observations), thus relying on communication to improve their performance can be inevitable. Furthermore, simulating the behaviour of agents of different preceptive abilities, that are operating under different interaction schemes, can have direct implications for the design of such agents by predicting their general performance and consequently their usefulness over different environmental settings.

This chapter is organised as follows. I will begin in Section 4.3 by briefly re-introducing my main approach towards measuring the performance of artificial agents. I also present in the same section a simplified measure that is used to abstract the amount of information, received by an agent as part of its observations or, transmitted by this agent via a communication data segment. I proceed in Section 4.5 by describing the details of my experiments whose purpose is to compare the influence of observation and communication on the performance of the evaluated agents. The outcomes from these experiments are then discussed in detail in Section 4.6. I conclude in Section 4.7 with a summary of the outcomes from this chapter and some of its implications on the current state of research.

4.3 Measuring Information

I will use the Λ^* environment described in Section 3.6 to conduct a series of controlled experiments on a group of interactive agents operating in this environment. The principal idea is to adjust the evaluated agents' communication and observation ranges and record whether (and how much) these changes impact the agents' measured intelligence.

I will continue to follow the agent-environment framework (Legg and Hutter, 2007) described in Section 3.5 to run my experiments. In this setting the environment is the world where agents can interact using a set of observations, actions and rewards. Recall that, at each step or iteration of the test, the environment generates observations from the set of observations \mathcal{O} and sends them to the agents. Agents then perform actions from a limited set of actions \mathcal{A} in response. Finally, the environment rewards back each agent from the set $\mathcal{R} \subseteq \mathbb{Q}$ based on the quality of its action. Thus, each iteration or step of the test stands for one sequence of observation-action-reward.

Shannon's *entropy* (Shannon, 1948) is once more used to measure the *uncertainty* $\mathcal{H}(\mu)$ in a given environment μ (recall Section 3.6.3). However, for simplicity, we redefine the set of possible states N of an (instance of the Λ^*) environment μ as the set of all possible cell positions that special object \oplus can be in. Therefore, in an m -by- n environment space $N = mn$. At the beginning of a test, the entropy is maximal as there is complete uncertainty about the current state of μ from an agent's perspective. Therefore the probability $p(s_\mu)$ of a given state s_μ occurring follows a uniform distribution and is equal to $1/|N|$. Using \log_2 as a base for calculations, the uncertainty $\mathcal{H}(\mu)$ is calculated as $\mathcal{H}(\mu) = - \sum_{s_\mu \in N} p(s_\mu) \log_2 p(s_\mu) = \log_2 |N|$ bits.

In addition to the above, Shannon's entropy will also be used in this chapter to abstract the amount of information received in an observation o from the environment, as well as to abstract the amount of information transmitted by an agent through a communication message c .

The environment sends observations to the agents holding a description of (the rewards in) their range of *k-Moore neighbourhood* (Gray, 2003; Weisstein, 2015) cells, where $k \in \mathbb{N}$. The number of cells contained in a range of k -Moore neighbourhood denoted by $C_k^{Moore} = (2k + 1)^2$ cells. Therefore, the larger the observation range of an agent, the more information it is given about the state of the environment. For simplicity¹⁷, I abstract/measure the amount of information an agent π is given about the environment as the entropy $H(o)$ of the observation o sent to π by the environment at one iteration of the test, which translates to $\log_2 C_k^{Moore}$ bits, where C_k^{Moore} corresponds to the number of (neighbour) cells described in its observation o . Consequently, we expect from the theory that the more information given to an agent, the more likely that this agent will accurately reason about it by processing and interpreting the provided information.

¹⁷Note that this is a simplification using an abstraction of the amount of information. The actual amount of information can be measured as minimum the number of bits used to describe an observation o (or a communication message c) or as the logarithm of the probability distribution of the rewards in the environment after receiving o (or c). Therefore, different, and more elaborate, types of observation and communication data can be quantitatively captured and compared.

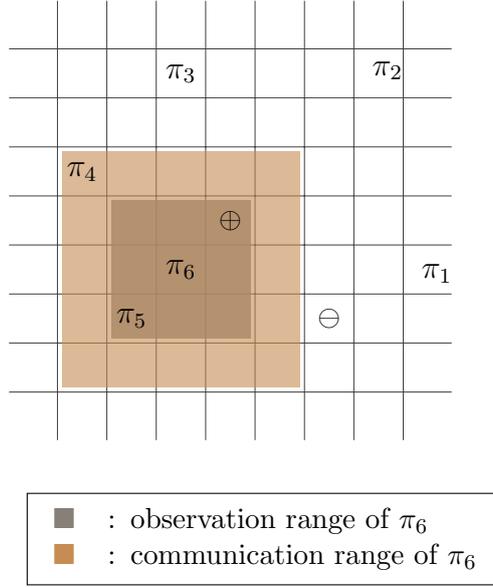


Figure 4.1: A illustration of a 9-by-9 Λ^* environment space of uncertainty $\mathcal{H}(\mu) = \log_2 81 = 6.33$ bits. Agent π_6 has an observation range of $k = 1$ Moore neighbourhood and a communication range of $k' = 2$ Moore neighbourhood depicted as the grey and brown shaded areas respectively. The observation and communication entropies are thus calculated as $H(o) = \log_2(2 \times 1 + 1)^2 = 3.16$ bits and $H(c) = \log_2(2 \times 2 + 1)^2 = 4.64$ bits respectively. In this figure, agent π_6 can only communicate with agents π_4 and π_5 .

Likewise, the amount of information transmitted by an agent π in a communication message c is calculated as the entropy $H(c)$ which is in turn calculated as $\log_2 C_{k'}^{Moore}$ bits, where $C_{k'}^{Moore}$ corresponds here to the number of cells a communication message c is transmitted over. Hence, the larger the communication range of an agent the more data it can transmit over that range. An example depicting the communication and observation ranges of an agent operating in the Λ^* environment is given in Figure 4.1.

In the context of the above definitions, and the agent-environment framework (Section 3.5), a group of interactive agents can be evaluated over a Λ^* environment μ (at one test iteration) as follows:

1. The environment μ sends an observation o to each agent in the group.
2. Each agent communicates or interacts with other agents located within its communication range by transmitting (and receiving from other agents) some data c over all the cells located within this range.
3. Each agent then performs an action based on the details received from its observation and communication data.
4. The environment rewards back each agent according to the quality of its action.

In order to compare the influence of the factors of observation and communication on performance, the observation and communication entropies $H(o)$ and $H(c)$ can be altered by extending/narrowing the observation and communication ranges of the Moore neighbourhoods respectively. This will enable us to conduct new sorts of experiments that are

different from the ones conducted in the previous chapter in which we exclusively tested over a range of $k = 1$ Moore neighbourhood.

4.4 Agent Interaction Modes

For the purpose of this chapter, I will evaluate groups of local search agents using a variation of three of the communication (or interaction) techniques defined in Section 3.9 of the previous chapter, namely (i) stigmergy or indirect communication (Section 3.9.1), (ii) implicit leadership through auctions and bidding (Section 3.9.2) and (iii) imitating a super-solver or expert agent (Section 3.9.3).

In the experiments conducted in this chapter, the agents communicating using implicit leadership through auctions and bidding (Section 3.9.2) will be restricted to bid against other agents located in their communication range only as opposed to all other agents operating in the same environment. For brevity I will refer to this communication mode as *direct communication* for the rest of this chapter since the agents using this communication technique bid visibly to one another. Note that the (individual) local search agent behaviour defined in Section 3.8.1 holds and will not be repeated here. As for the details of the communication modes, an informal and concise description is given below.

- (i) Stigmergy or indirect communication: agents coordinate indirectly by altering the environment so that it reflects their individual observations. At each iteration of the test, when an agent senses a reward as part of its observation o , it puts *fake-rewards* in all the cells within its communication range. The fake rewards reflect the real reward the agent has observed.
- (ii) Direct communication: at each iteration of the test, agents broadcast a description of their highest observed reward to all other agents in their communication range. The agent with the highest visible reward wins the auction and selects the target cell to be approached by all its neighbouring agents located in its communication range. All other agents (that are not located within the communication range of the winner) perform according to the individual local search behaviour.
- (iii) Imitation: in this setting (in addition to the evaluated agents) I introduce an un-evaluated expert agent that always takes the most rewarding action at each iteration of the test. The evaluated local search agents imitate the expert agent by mimicking its action when it happens to be located within their communication range. Each agent in turn shares (or propagates) its observation with other agents located within its communication range, if any exist. Local search agents that are not located within the communication range of the expert agent follow the individual local search behaviour.

4.5 Experimental Setup

I conducted a series of controlled experiments where I evaluate a collective of 20 cooperative local search agents, denoted by Π , over a Λ^* environment μ such that $\mathcal{H}(\mu) = 11.28$ bits of uncertainty. Each experiment consisted of 200 iterations of *observation-communication-action-reward* sequences.

The experiments are repeated across a range of entropy values for $H(o)$ and $H(c)$ between 0.04 and 10.84 bits. As described earlier, this was done by changing the observation and communication ranges of the agents. I recorded the scores of the collective, $\mathcal{Y}(H_o, H_c) \in [-1.0, 1.0]$, over all the different values of $H(o)$ and $H(c)$ used in my experiments. I denote by E the set of entropy values of $H(o)$ and $H(c)$ that were used in these experiments. The average scores of collective Π from 1000 repeated experiments¹⁸ are depicted in figures 4.2, 4.5 and 4.8. Each one of these figures shows the scores of collective Π using one of the three communication modes described earlier in Section 4.4, across all the entropy values in E . For clarity I use the notation Π^i, Π^{ii} and Π^{iii} to refer to collective Π when using communication modes (i), (ii) and (iii) from Section 4.4 respectively.

For alternative visualisations that allow easier and more thorough interpretation of these results, I have re-plotted the scores $\mathcal{Y}(H_o, H_c)$ illustrated in figures 4.2, 4.5 and 4.8 for fixed values of $H(c)$ across increasing values of $H(o)$ and vice versa. The resulting plots (among others) are shown in Figures 4.3, 4.6 and 4.9 respectively. I analyse and discuss the results from the above experiments in the next section.

4.6 Results and Discussion

The scores of collectives Π^i, Π^{ii} and Π^{iii} using (i) stigmergy or indirect communication, (ii) direct communication and (iii) imitation are discussed in Sections 4.6.1, 4.6.2 and 4.6.3 all respectively.

4.6.1 Indirect communication

Figures 4.2 and 4.3 show that the effectiveness of the agents monotonically increases along with their observation entropy $H(o)$ until it converges around an $H(o)$ of 10.8 bits.

Increasing the (stigmergic) communication entropy $H(c)$ between the agents also has an impact on their intelligence. However, the influence of $H(c)$ on intelligence is rather more complicated as it seems to also depend on the observation entropies $H(o)$. For instance, for an $H(o)$ of 0.04 bits, the best performance, $\max(\mathcal{Y}(H_o, H_c))$, is reached when the coefficient $\alpha = \frac{H(c)}{H(o)} = 9$. For larger $H(o)$ entropies, the best performances are reached at smaller α values, and $\alpha \rightarrow 1.0$ at an $H(o)$ of 10.84 bits. The overall picture from Figure 4.2 shows

¹⁸Note that the *coefficient of variation*—also known as the *relative standard deviation* (calculated as the ratio of standard deviation σ to the mean μ)—of the scores is less than 0.025 across the experiments, thus showing a low measure of variability of the score data between different experiments. I consider the small variations in scores along the fourth decimal place as experimental error.

		Observation entropy $H(o)$																							
		0.04	0.11	0.22	0.36	0.54	0.76	1.01	1.30	1.62	1.99	2.38	2.82	3.29	3.79	4.33	4.91	5.53	6.18	6.86	7.58	8.34	9.14	9.97	10.84
Communication entropy $H(c)$	10.84	0.0506	0.0693	0.0927	0.1120	0.1329	0.1517	0.1735	0.1998	0.2256	0.2545	0.2753	0.2903	0.3059	0.3277	0.3431	0.3615	0.3833	0.4076	0.4234	0.4328	0.4439	0.4470	0.4472	0.4470
	9.97	0.0469	0.0646	0.0837	0.1078	0.1290	0.1502	0.1712	0.1978	0.2256	0.2667	0.2990	0.3179	0.3297	0.3422	0.3609	0.3778	0.3933	0.4140	0.4291	0.4385	0.4448	0.4475	0.4468	0.4473
	9.14	0.0450	0.0657	0.0883	0.1008	0.1203	0.1434	0.1738	0.1950	0.2486	0.2907	0.3245	0.3419	0.3491	0.3652	0.3823	0.3930	0.4073	0.4213	0.4335	0.4398	0.4470	0.4472	0.4469	0.4465
	8.34	0.0441	0.0618	0.0824	0.1022	0.1200	0.1409	0.1766	0.2188	0.2637	0.3064	0.3394	0.3588	0.3695	0.3779	0.3955	0.4059	0.4180	0.4273	0.4339	0.4458	0.4470	0.4459	0.4464	0.4468
	7.58	0.0461	0.0702	0.0849	0.1046	0.1222	0.1497	0.1841	0.2333	0.2793	0.3249	0.3541	0.3707	0.3845	0.3945	0.4046	0.4131	0.4231	0.4311	0.4419	0.4450	0.4449	0.4466	0.4461	0.4471
	6.86	0.0489	0.0677	0.0863	0.1056	0.1274	0.1584	0.2024	0.2542	0.3105	0.3445	0.3696	0.3828	0.3943	0.4005	0.4111	0.4202	0.4268	0.4389	0.4410	0.4448	0.4458	0.4465	0.4463	0.4467
	6.18	0.0486	0.0671	0.0875	0.1076	0.1313	0.1676	0.2249	0.2817	0.3274	0.3575	0.3752	0.3897	0.4002	0.4079	0.4153	0.4229	0.4325	0.4383	0.4388	0.4409	0.4432	0.4438	0.4451	0.4468
	5.53	0.0478	0.0720	0.0927	0.1139	0.1470	0.1853	0.2503	0.3062	0.3457	0.3669	0.3827	0.3969	0.4053	0.4119	0.4172	0.4284	0.4327	0.4318	0.4360	0.4362	0.4386	0.4404	0.4438	0.4467
	4.91	0.0506	0.0721	0.0959	0.1199	0.1556	0.2088	0.2733	0.3255	0.3614	0.3770	0.3872	0.4016	0.4080	0.4141	0.4236	0.4295	0.4293	0.4306	0.4336	0.4353	0.4370	0.4408	0.4442	0.4468
	4.33	0.0506	0.0701	0.0965	0.1268	0.1726	0.2388	0.3002	0.3446	0.3698	0.3875	0.3941	0.4053	0.4136	0.4202	0.4246	0.4255	0.4271	0.4294	0.4299	0.4300	0.4330	0.4395	0.4436	0.4471
	3.79	0.0498	0.0745	0.1018	0.1413	0.1971	0.2685	0.3224	0.3579	0.3803	0.3893	0.3989	0.4073	0.4157	0.4187	0.4167	0.4197	0.4231	0.4227	0.4220	0.4249	0.4294	0.4373	0.4432	0.4469
	3.29	0.0528	0.0741	0.1050	0.1534	0.2286	0.2955	0.3363	0.3670	0.3859	0.3960	0.4037	0.4103	0.4165	0.4138	0.4168	0.4179	0.4162	0.4171	0.4158	0.4182	0.4256	0.4344	0.4425	0.4470
	2.82	0.0486	0.0790	0.1126	0.1686	0.2504	0.3151	0.3456	0.3732	0.3870	0.3982	0.4038	0.4132	0.4081	0.4109	0.4119	0.4135	0.4108	0.4071	0.4122	0.4128	0.4186	0.4308	0.4414	0.4464
	2.38	0.0533	0.0809	0.1203	0.1881	0.2643	0.3216	0.3506	0.3770	0.3906	0.4011	0.4066	0.4046	0.4056	0.4054	0.4018	0.3999	0.4014	0.3987	0.4034	0.4102	0.4165	0.4279	0.4403	0.4463
	1.99	0.0530	0.0836	0.1338	0.2069	0.2717	0.3202	0.3493	0.3731	0.3897	0.3976	0.3947	0.4017	0.4001	0.3958	0.3948	0.3917	0.3913	0.3956	0.3948	0.4023	0.4113	0.4285	0.4404	0.4465
1.62	0.0533	0.0898	0.1457	0.2127	0.2681	0.3125	0.3444	0.3705	0.3893	0.3902	0.3907	0.3919	0.3939	0.3867	0.3846	0.3856	0.3862	0.3843	0.3869	0.3990	0.4083	0.4269	0.4404	0.4461	
1.30	0.0570	0.0929	0.1562	0.2140	0.2566	0.2991	0.3378	0.3690	0.3729	0.3830	0.3816	0.3879	0.3829	0.3784	0.3774	0.3800	0.3791	0.3811	0.3873	0.3899	0.4062	0.4238	0.4391	0.4469	
1.01	0.0565	0.1033	0.1582	0.2100	0.2538	0.2894	0.3195	0.3447	0.3622	0.3741	0.3772	0.3811	0.3755	0.3741	0.3739	0.3752	0.3779	0.3805	0.3830	0.3926	0.4056	0.4247	0.4397	0.4468	
0.76	0.0604	0.1051	0.1568	0.1969	0.2339	0.2635	0.2985	0.3269	0.3488	0.3615	0.3665	0.3693	0.3705	0.3718	0.3727	0.3745	0.3757	0.3811	0.3799	0.3886	0.4017	0.4243	0.4403	0.4467	
0.54	0.0645	0.1002	0.1448	0.1849	0.2159	0.2478	0.2758	0.3061	0.3352	0.3500	0.3577	0.3635	0.3705	0.3736	0.3757	0.3771	0.3747	0.3794	0.3813	0.3890	0.4008	0.4220	0.4392	0.4468	
0.36	0.0660	0.0956	0.1355	0.1637	0.1899	0.2197	0.2580	0.2882	0.3188	0.3381	0.3560	0.3650	0.3705	0.3741	0.3781	0.3798	0.3828	0.3873	0.3853	0.3899	0.4022	0.4216	0.4405	0.4467	
0.22	0.0605	0.0917	0.1150	0.1514	0.1777	0.2061	0.2378	0.2714	0.3023	0.3317	0.3498	0.3643	0.3768	0.3797	0.3817	0.3877	0.3910	0.3916	0.3941	0.3958	0.4058	0.4231	0.4420	0.4472	
0.11	0.0548	0.0779	0.0997	0.1313	0.1592	0.1893	0.2186	0.2520	0.2864	0.3165	0.3372	0.3630	0.3728	0.3845	0.3911	0.3952	0.4016	0.4032	0.4024	0.4090	0.4137	0.4274	0.4422	0.4469	
0.04	0.0486	0.0676	0.0919	0.1146	0.1400	0.1715	0.2025	0.2317	0.2671	0.2987	0.3307	0.3542	0.3734	0.3890	0.3970	0.4031	0.4091	0.4116	0.4141	0.4194	0.4265	0.4305	0.4422	0.4464	

Figure 4.2: A plot of the test scores $\Upsilon(H_o, H_c)$ from the experiments described in Section 4.5. Experiments are conducted across different $H(o)$ and $H(c) \in E$ values (in bits) for collective Π^i using **indirect communication by stigmergy** (corresponding to communication mode (i) in Section 4.4). The grey colour-map intensities reflect how large the score values $\Upsilon(H_o, H_c)$ are (in each table cell/experiment). For instance, strong colour intensities mean larger scores. Large score values are shaded with dark grey or black whereas low scores are shaded with lighter colours (Chmait, Dowe, Green and Li, 2015, Fig. 1).

that the performance drops as the entropy $H(c)$ moves away from $\alpha \times H(o)$. This non-monotonic variation of scores shows that increasing communication does not necessarily always lead to an increase in performance as generally presumed.

To better understand the influence of indirect communication on the scores of the collectives we have to analyse further the relationship between $H(o)$ and $H(c)$. Figure 4.4 is a whisker plot showing the variation in the scores (depicted in Figure 4.2) across different entropy values $H(c) \subseteq E$ for fixed entropies $H(o)$, and vice versa. The central mark (shown in red in Figure 4.4) is the median while the edges of the box represent the 25th and 75th percentiles of the scores and the whiskers extend to the most extreme score values. The blue line-plot shows the average scores at each of the intermediate entropy values.

Figure 4.4 shows that—for indirect communication— $H(c)$ is most significant when $H(o) \in [0.3, 1.9]$ bits. For instance, using stigmergy to communicate very short observations (low $H(o)$ entropies) does not have a large influence on performance, possibly because the observations do not carry much information. Likewise, using stigmergy within collectives of agents with extended observation abilities (high $H(o)$ entropies) has no significant effect on performance, as the uncertainty in the environment has already significantly been reduced as a result of the agents' broad observations. However, communication using stigmergy was fairly effective in less extreme cases. To make my observations more concrete, I define below the communication-over-observation coefficient of success ϕ .

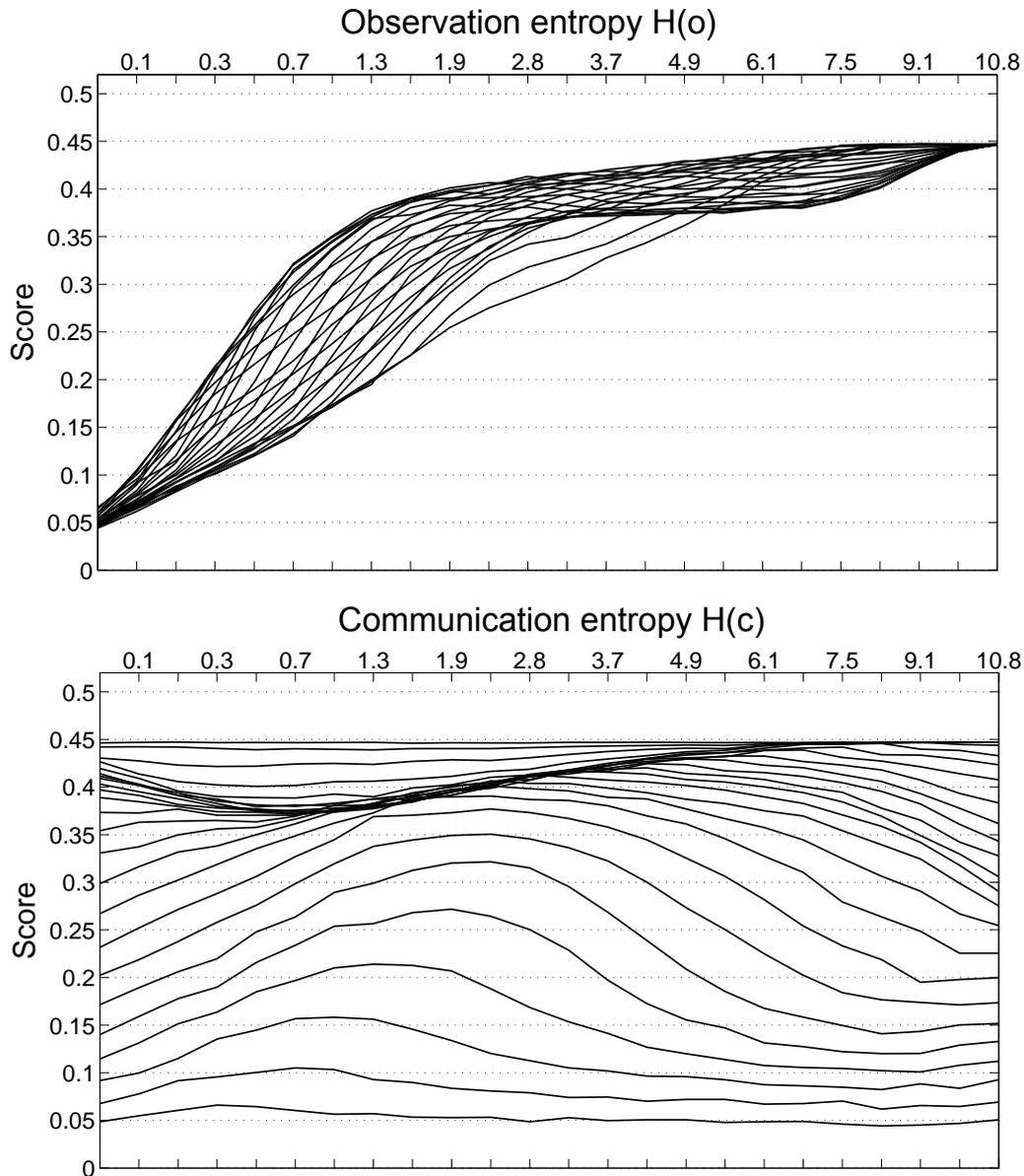


Figure 4.3: Variation in the scores (from Figure 4.2) of collective Π^i using indirect communication (stigmergy) (Chmait, Dowe, Green and Li, 2015, Fig. 2). The top plot shows the scores for fixed values of $H(c)$ across increasing values of $H(o)$, that is for each row of the entries in the table in Figure 4.2. The bottom plot shows the scores for fixed values of $H(o)$ across increasing values of $H(c)$, that is for each column of the entries in the table in Figure 4.2.

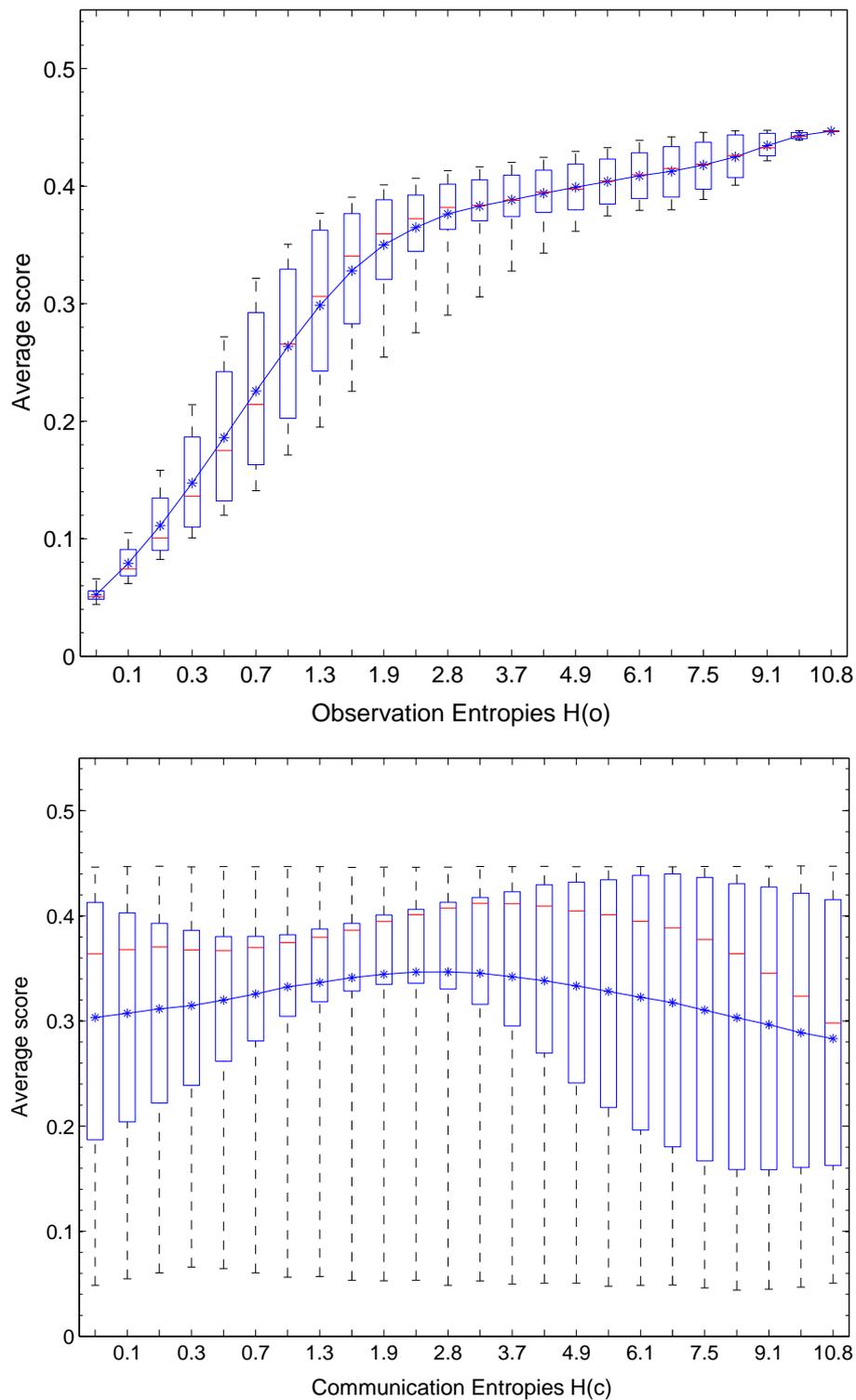


Figure 4.4: Whisker plot showing the variation in the test scores of agents relying on indirect communication from Figure 4.2 (or equivalently Figure 4.3). The whiskers show the variation in scores across different entropy values $H(c)$, for fixed entropies $H(o)$ (top plot), and vice versa (bottom plot). The central mark (shown in red) is the median while the edges of the box represent the 25th and 75th percentiles of the scores and the whiskers extend to the most extreme score values. The blue line-plot shows the average scores at each of the intermediate entropy values (Chmait, Dowe, Green and Li, 2015, Fig. 3).

Definition 4 Let $S = \{(x, y) \in E \times E \mid x > y\}$. The communication-over-observation coefficient of success is:

$$\phi = \left(\sum_S \text{inferior}(\mathcal{Y}(x, y), \mathcal{Y}(y, x)) \div |S| \right)$$

where $\text{inferior}(a, b)$ is a function that returns 1 if $a < b$, or zero otherwise.

For this mode of communication, the coefficient $\phi = 11/276 = 0.0399$. By taking into account that the test scores are of the form $\mathcal{Y}(H_o, H_c)$, the value of ϕ suggests that, for this communication mode, it is much more effective to increase the observation entropies of the agents as opposed to increasing their communication entropies¹⁹. More importantly, the dependency of communication $H(c)$ on observation $H(o)$ is made explicit here. For instance, using $H(c)$ values inferior to $H(o)$ is rarely more rewarding than in the reciprocal case.

4.6.2 Direct communication

In this section I perform similar analysis as above on the scores of collective Π^{ii} implementing a direct communication mode between its members. The scores for this collective are given in Figure 4.6. By comparing Figures 4.3 and 4.6, we observe that, while increasing observation entropies still leads to a significant increase in performance, the influence of direct communication is much more significant than in the case of indirect communication. We can see a clear pattern in Figure 4.6 showing higher performances for higher communication entropies for a fixed $H(o)$.

We also observe that, in this setting, using very low $H(o)$ entropies does not ensure an optimal performance for Π^{ii} . However, re-compensating the short-sighted agents of Π^{ii} with high $H(c)$ entropies can lead to a group that is up to four times better in performance. For instance, introducing communication in the collective has allowed the scores to increase from around 0.1 to roughly 0.45 as illustrated in the bottom plot of Figure 4.6. This also indicates the very low-dependency of $H(c)$ on the value of $H(o)$.

On the other hand, for fairly high observation entropies, augmenting communication between the agents is at least as effective as mounting their observations, and can sometimes be even more effective as shown in Figure 4.7. In this setting the coefficient $\phi = 272/276 = 0.9855$ (recall Definition 4), meaning that augmenting the communication entropies within the system will highly likely lead to a more intelligent system. Consequently, communication is effective here even when the observation entropies are slim, again suggesting a low dependency on $H(o)$.

¹⁹Recall that we are experimenting for the entropy values E , using a number of agents $|\Pi| = 20$, over an environment of uncertainty $H(\mu) = 11.28$ bits. Experimenting with different number of agents and/or environment spaces can result in the measured (optimal) performances to occur at different entropy values, although the general picture is the same.

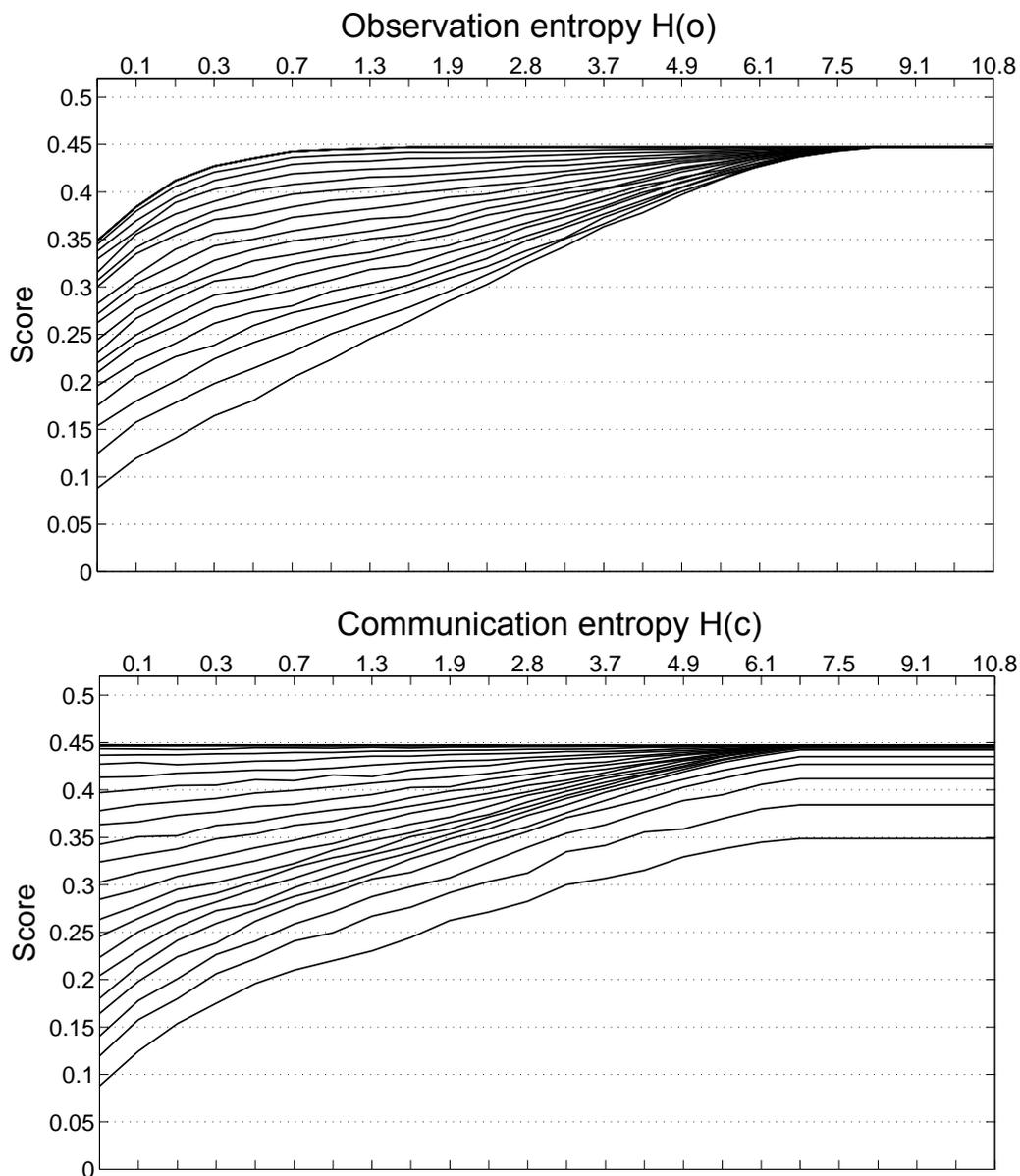


Figure 4.6: Variation in scores for collective Π^{ii} using the direct communication mode (Chmait, Dowe, Green and Li, 2015, Fig. 2). The top plot shows the scores for fixed values of $H(c)$ across increasing values of $H(o)$, that is for each row of the entries in the table in Figure 4.5. The bottom plot shows the scores for fixed values of $H(o)$ across increasing values of $H(c)$, that is for each column of the entries in the table in Figure 4.5.

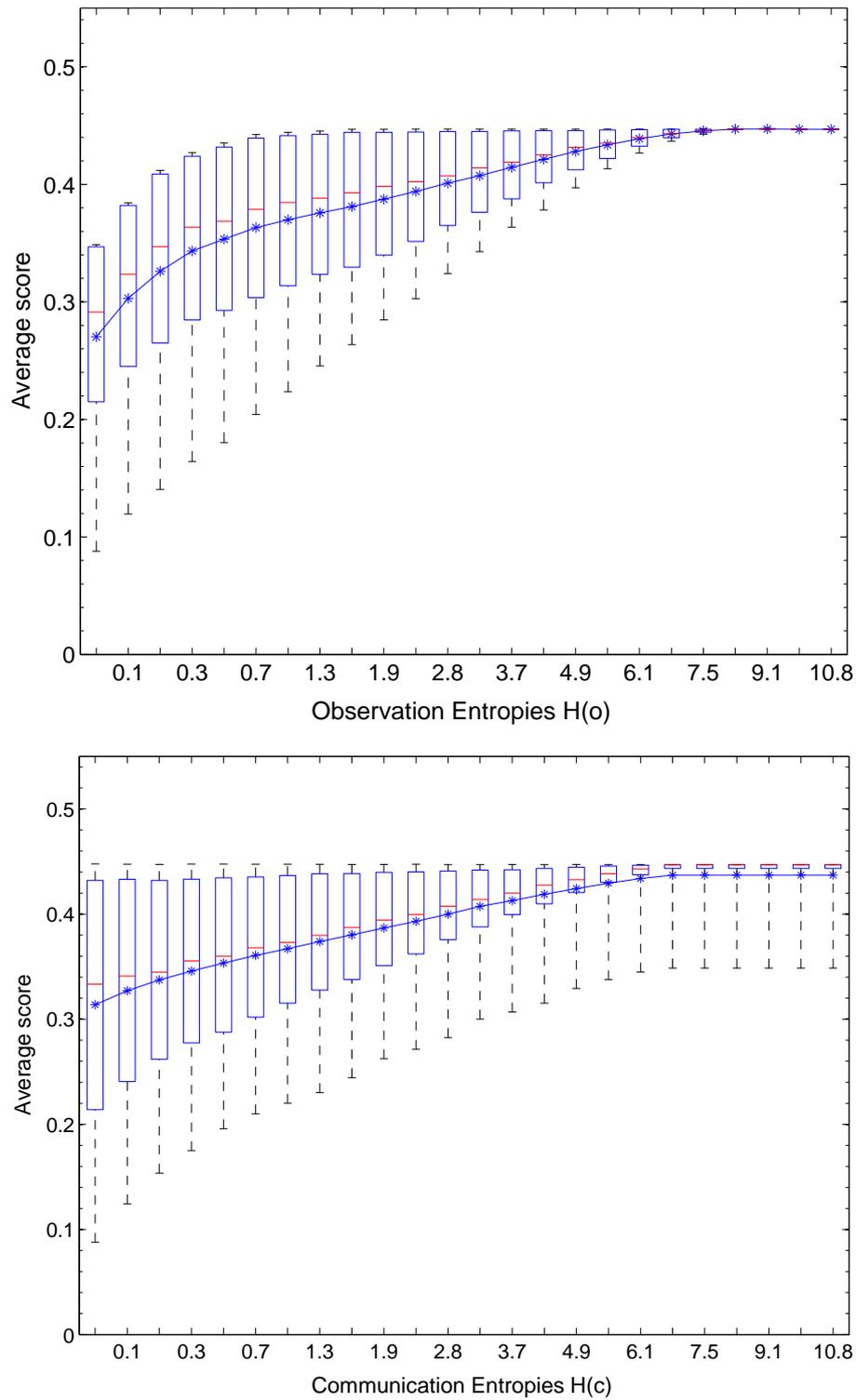


Figure 4.7: Whisker plot showing the variation in the test scores of collective Π^{ii} relying on direct communication taken from Figure 4.5 (or equivalently Figure 4.6). The whiskers show the variation in scores across different entropy values $H(c)$, for fixed entropies $H(o)$ (top plot), and vice versa (bottom plot) (Chmait, Dowe, Green and Li, 2015, Fig. 3).

		Observation entropy H(o)																								
		0.04	0.11	0.22	0.36	0.54	0.76	1.01	1.30	1.62	1.99	2.38	2.82	3.29	3.79	4.33	4.91	5.53	6.18	6.86	7.58	8.34	9.14	9.97	10.84	
Communication entropy H(c)	10.84	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	
	9.97	0.7977	0.7977	0.7977	0.7978	0.7977	0.7978	0.7978	0.7978	0.7978	0.7977	0.7978	0.7978	0.7977	0.7977	0.7977	0.7977	0.7978	0.7978	0.7978	0.7978	0.7977	0.7977	0.7978	0.7978	
	9.14	0.7976	0.7976	0.7976	0.7977	0.7976	0.7977	0.7976	0.7976	0.7976	0.7977	0.7976	0.7976	0.7977	0.7977	0.7977	0.7977	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	
	8.34	0.7975	0.7975	0.7975	0.7976	0.7976	0.7976	0.7976	0.7976	0.7976	0.7976	0.7976	0.7976	0.7977	0.7977	0.7977	0.7977	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	0.7978	
	7.58	0.7977	0.7975	0.7977	0.7975	0.7977	0.7974	0.7976	0.7976	0.7976	0.7975	0.7975	0.7977	0.7976	0.7977	0.7978	0.7978	0.7978	0.7979	0.7977	0.7979	0.7980	0.7978	0.7978	0.7978	
	6.86	0.7973	0.7976	0.7975	0.7974	0.7973	0.7973	0.7976	0.7975	0.7977	0.7976	0.7975	0.7976	0.7973	0.7976	0.7974	0.7978	0.7979	0.7979	0.7979	0.7979	0.7980	0.7979	0.7979	0.7978	
	6.18	0.7973	0.7971	0.7973	0.7969	0.7972	0.7971	0.7976	0.7973	0.7973	0.7974	0.7971	0.7973	0.7972	0.7976	0.7976	0.7976	0.7976	0.7976	0.7976	0.7977	0.7980	0.7981	0.7978	0.7979	
	5.53	0.7973	0.7972	0.7971	0.7976	0.7973	0.7971	0.7974	0.7974	0.7974	0.7975	0.7972	0.7975	0.7974	0.7972	0.7972	0.7979	0.7979	0.7979	0.7979	0.7979	0.7979	0.7979	0.7978	0.7978	
	4.91	0.7966	0.7964	0.7966	0.7966	0.7965	0.7968	0.7970	0.7972	0.7966	0.7970	0.7969	0.7967	0.7970	0.7975	0.7974	0.7974	0.7974	0.7978	0.7979	0.7976	0.7980	0.7978	0.7978	0.7978	
	4.33	0.7958	0.7961	0.7960	0.7960	0.7960	0.7962	0.7964	0.7964	0.7966	0.7968	0.7968	0.7968	0.7974	0.7975	0.7969	0.7978	0.7977	0.7977	0.7983	0.7981	0.7980	0.7982	0.7981	0.7981	
	3.79	0.7950	0.7956	0.7953	0.7956	0.7954	0.7955	0.7958	0.7955	0.7958	0.7963	0.7967	0.7960	0.7959	0.7972	0.7967	0.7970	0.7969	0.7970	0.7973	0.7976	0.7979	0.7982	0.7977	0.7979	
	3.29	0.7937	0.7952	0.7945	0.7947	0.7951	0.7938	0.7947	0.7949	0.7949	0.7957	0.7952	0.7949	0.7955	0.7966	0.7960	0.7962	0.7968	0.7972	0.7971	0.7979	0.7981	0.7979	0.7979	0.7979	
	2.82	0.7916	0.7913	0.7915	0.7931	0.7916	0.7922	0.7925	0.7934	0.7922	0.7944	0.7929	0.7943	0.7938	0.7955	0.7953	0.7954	0.7962	0.7962	0.7973	0.7970	0.7988	0.7976	0.7975	0.7976	
	2.38	0.7851	0.7834	0.7821	0.7846	0.7833	0.7856	0.7859	0.7849	0.7837	0.7852	0.7871	0.7895	0.7900	0.7900	0.7914	0.7937	0.7932	0.7948	0.7968	0.7971	0.7983	0.7976	0.7972	0.7975	
	1.99	0.7610	0.7593	0.7615	0.7633	0.7619	0.7593	0.7587	0.7654	0.7641	0.7618	0.7702	0.7655	0.7721	0.7737	0.7793	0.7815	0.7882	0.7918	0.7950	0.7958	0.7972	0.7977	0.7968	0.7970	
	1.62	0.7127	0.7040	0.7115	0.7101	0.7045	0.7092	0.7086	0.7039	0.7062	0.7140	0.7183	0.7253	0.7337	0.7497	0.7555	0.7615	0.7754	0.7809	0.7895	0.7913	0.7958	0.7967	0.7953	0.7952	
	1.30	0.6097	0.6166	0.6097	0.6138	0.6055	0.6237	0.6146	0.6345	0.6267	0.6436	0.6527	0.6614	0.6854	0.7006	0.7162	0.7322	0.7526	0.7682	0.7812	0.7874	0.7930	0.7925	0.7917	0.7921	
	1.01	0.5052	0.5024	0.5071	0.5055	0.5126	0.5064	0.5218	0.5310	0.5573	0.5743	0.5996	0.6151	0.6418	0.6680	0.6911	0.7166	0.7427	0.7603	0.7730	0.7844	0.7891	0.7912	0.7902	0.7898	
	0.76	0.4050	0.4179	0.4161	0.4078	0.4130	0.4397	0.4653	0.4780	0.5062	0.5369	0.5583	0.5941	0.6191	0.6506	0.6790	0.7054	0.7311	0.7566	0.7704	0.7792	0.7867	0.7874	0.7867	0.7867	
	0.54	0.3244	0.3319	0.3414	0.3232	0.3570	0.3796	0.4192	0.4419	0.4697	0.5038	0.5380	0.5718	0.6015	0.6442	0.6649	0.6982	0.7248	0.7471	0.7634	0.7771	0.7829	0.7837	0.7818	0.7815	
	0.36	0.2465	0.2639	0.2783	0.2940	0.3215	0.3555	0.3781	0.4181	0.4469	0.4841	0.5207	0.5547	0.5910	0.6277	0.6591	0.6891	0.7136	0.7397	0.7580	0.7699	0.7783	0.7774	0.7772	0.7771	
	0.22	0.1997	0.2101	0.2366	0.2667	0.2938	0.3299	0.3571	0.3970	0.4341	0.4675	0.5009	0.5394	0.5784	0.6125	0.6492	0.6783	0.7058	0.7336	0.7499	0.7657	0.7713	0.7726	0.7722	0.7725	
	0.11	0.1589	0.1932	0.2122	0.2377	0.2735	0.3148	0.3483	0.3819	0.4159	0.4624	0.4988	0.5302	0.5725	0.6052	0.6417	0.6744	0.7008	0.7270	0.7463	0.7599	0.7662	0.7676	0.7668	0.7671	
	0.04	0.1319	0.1672	0.2036	0.2318	0.2640	0.3031	0.3337	0.3740	0.4166	0.4458	0.4883	0.5272	0.5632	0.5998	0.6343	0.6666	0.6975	0.7198	0.7419	0.7544	0.7615	0.7625	0.7615	0.7620	Imitation

Figure 4.8: Test scores $\mathcal{T}(H_o, H_c)$ across different $H(o)$ and $H(c)$ (in bits) values for the agent collective Π^{iii} using imitation (Chmait, Dowe, Green and Li, 2015, Fig. 1). [The interpretation of colour intensities is the same as given in Figure 4.2.].

entropy values calculated according to Eq. (4.1) below:

$$\nabla \mathcal{T}(H_o, H_c) \leftarrow \left| \frac{\partial \mathcal{T}(H_o, H_c)}{\partial H(c)} \right| - \left| \frac{\partial \mathcal{T}(H_o, H_c)}{\partial H(o)} \right| \quad (4.1)$$

The outcome from Eq. (4.1) highlights the entropies where (in a environment of uncertainty $H(\mu) = 11.28$ bits, and using a number of agents $|\Pi| = 20$) communication has the highest influence on the effectiveness $\mathcal{T}(H_o, H_c)$ of Π when compared to the influence of observation. We observe that indirect communication records the highest impact across the entropies ranging in $[0.3, 1.9]$ bits. Direct communication is most significant within entropies of $[0.1, 1.9]$ bits, while imitation has the highest influence over entropy values in the range $[0.7, 4.9]$ bits²⁰. These types of results are relevant to the design of intelligent agent collectives. By fine-tuning the communication ranges of the agents with respect to their observation abilities one can improve the overall effectiveness of these agents and further measure how much advantageous can these agents approximately be.

4.7 Conclusion

This chapter is an extension to Chapter 3 in which I primarily tackled goal **(G01)**. In essence, I elaborate in this chapter on the third objective of this thesis, *Obj03*, which consists of analysing how different factors affect the performance of interactive agents, and how these factors are related to one another.

The main idea was to evaluate collaborative groups of artificial agents across different communication settings and observation abilities. By analysing the effectiveness of these

²⁰Recall footnote 19.

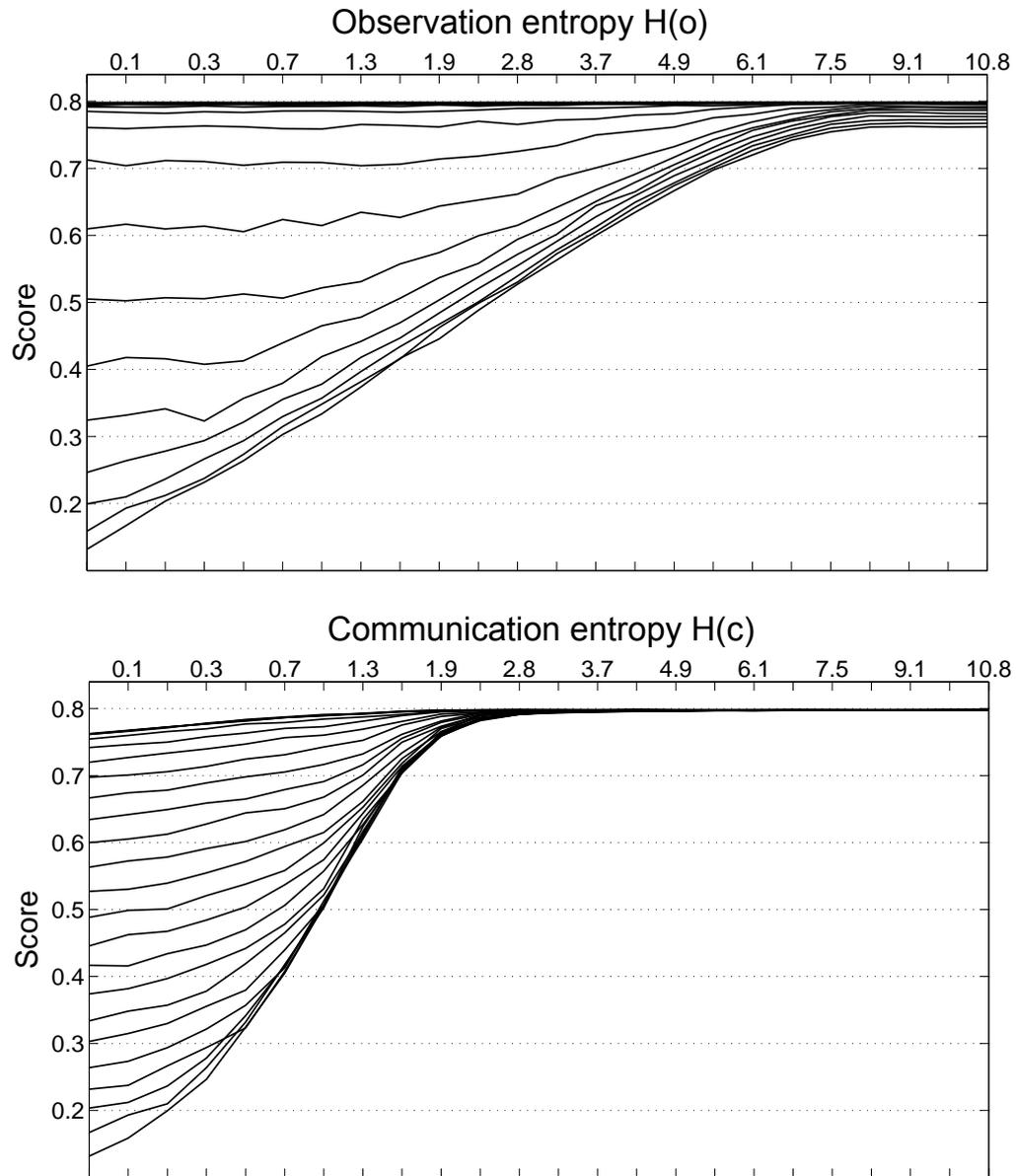


Figure 4.9: Variation in scores for collective Π^{iii} using imitation. The top plot shows the scores for fixed values of $H(c)$ across increasing values of $H(o)$, that is for each row of the entries in the table in Figure 4.8 (Chmait, Dowe, Green and Li, 2015, Fig. 2). The bottom plot shows the scores for fixed values of $H(o)$ across increasing values of $H(c)$, that is for each column of the entries in the table in Figure 4.8.

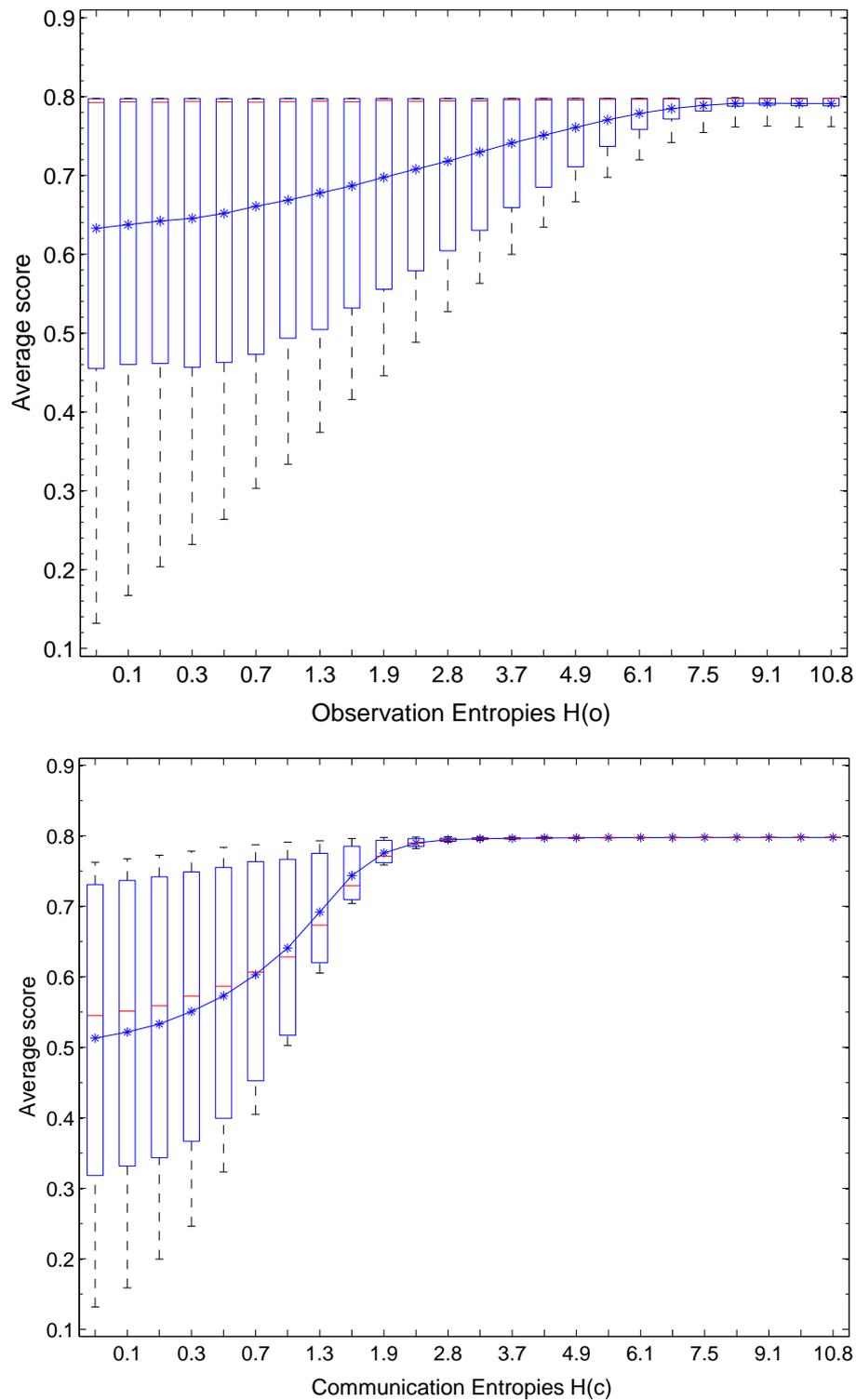


Figure 4.10: Whisker plot showing the variation in the test scores of collective Π^{iii} relying on imitation taken from Figure 4.8 (or equivalently Figure 4.9). The whiskers show the variation in scores across different entropy values $H(c)$, for fixed entropies $H(o)$ (top plot), and vice versa (bottom plot) (Chmait, Dowe, Green and Li, 2015, Fig. 3).

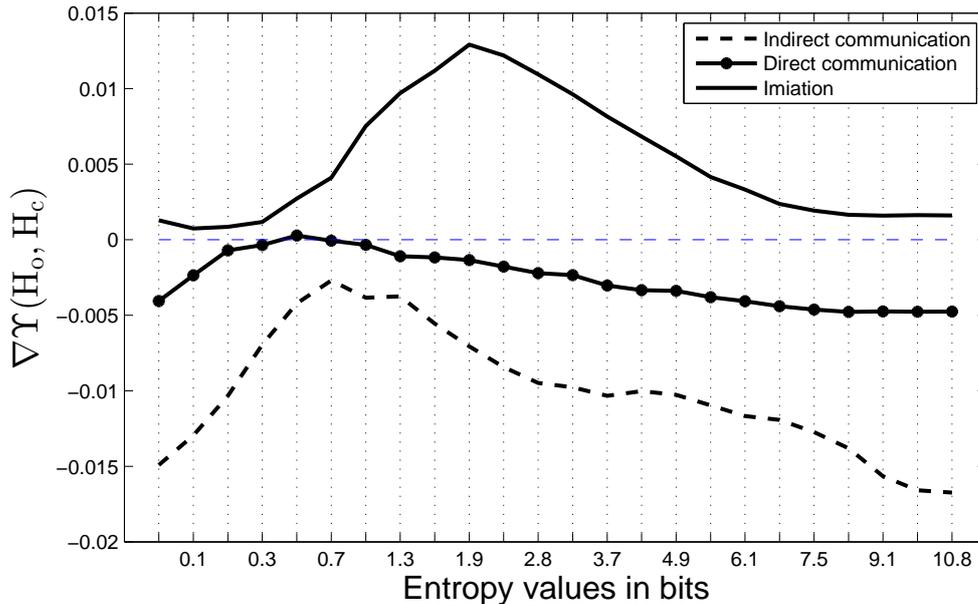


Figure 4.11: Average difference in gradient $\nabla\Upsilon(H_o, H_c)$ in $H(c)$ and $H(o)$ directions over a set of entropy values E for all three communication modes (Chmait, Dowe, Green and Li, 2015, Fig. 4).

groups, I measure the influence of two factors inherent to multiagent systems, namely the **observation and communication abilities** of their agents, on the overall intelligence of the evaluated system.

I highlighted the different configurations where the effectiveness of agent-based systems is significantly influenced by these two factors. I also discussed how dull systems with low observation or perception abilities can be re-compensated for, and significantly improved, by increasing the communication entropies of their agents, thus leading to smarter systems. Moreover, I identified circumstances where the increase in communication does not monotonically improve performance as commonly presumed.

In addition, the outcomes from my experiments gave rise to an important observation with respect to the relationship between an agent’s ability to perceive its environment, its ability to communicate and its performance. For instance, a group of agents implementing different types of communication modes between its members does not exhibit a uniform (overall) shift in performance over equivalent (communication and observation) entropy values. This reveals the existence of a certain type of dependency between communication and observation—which I measured and analysed.

Apart from improving our understanding of the observation and communication abilities of agents and their characteristics, the results from this chapter can have several implications for real world applications. For example the problem of exploring an unknown environment is frequently addressed in current research, e.g. (Laguna et al., 2014). Mobile robots are usually equipped with sensors and transmitters for communication and are dispatched to discover unknown or partially-known environments from large new planets and newly discovered territories to microscopic environments. A major challenge arises

when the robots' observation data does not provide complete coverage of the space of interest. Having a clear understanding of the robots' observation and communication abilities, and how they are related to one another, can make a big difference with respect to their collective performance.

Furthermore, the outcomes are helpful to understand some existing (biological) systems and their decision making. For instance, according to the results from my experiments, a simple stigmergy-like behaviour was (only significantly) effective when focused alterations to the environment took place (not too narrow or too broad with respect to the agents' observation range). This sort of behaviour can lead to a high concentration of agents repeatedly targeting the altered areas of their environment, in turn leading to a swarming behaviour. This is the case of many insect behaviours such as ants which end up travelling along the same path due to the accumulation of high pheromone levels on that path. However, I note that biological agents are not *synonymous* or completely equivalent to the types of artificial agents discussed here. For instance, biological agents are constrained in perception (local, diverse and disturbed), cognition (limited, heterogeneous and biased) and communication (limited, unreliable, incomplete, deceptive and strategic), integration and action selection, which is not always the case in my experiments. In my future work, I intend to explore more communication protocols reflecting bottom-up sensing and top-down control reflecting the flow of information in biological systems and organisations.

The end of chapter 4 serves as a closure to my first goal (**G01**) and in turn its three enclosed objectives *Obj01*, *Obj02* and *Obj03*. In the next chapter I will use a different approach in order to mathematically model the accuracy of agent-based systems.

Chapter 5

A Mathematical Model for Predicting the Accuracy of AI Agents

Mathematical reasoning may be regarded rather schematically as the exercise of a combination of two faculties, which we may call intuition and ingenuity.

—Alan Turing, Systems of logic based on ordinals (1938)

The research in this chapter has been published in the following article:

- Nader Chmait, David L. Dowe, David, Yuan-Fang Li and G. Green (2017). *An Information Theoretic Predictive Model for the Accuracy of AI Agents Adopted From Psychometrics*, Proceedings of the 10th International Conference on Artificial General Intelligence (AGI), Vol. 10414 of Lecture Notes in Artificial Intelligence (LNAI), Chapter 21, Melbourne, Australia, Springer. [**Winner of the 2017 Kurzweil Best Paper Prize**]. https://doi.org/10.1007/978-3-319-63703-7_21

5.1 Overview

Intelligence tests are powerful tools for assessing subjects of interest of different abilities and comparing their performances. Despite their popularity, the outcomes from intelligence tests are limited to generating an overall measure of accuracy, usually across a range of task difficulties. That is to say, while intelligence scores are accurate measures of some abilities associated with the evaluated agent, they are an unreliable predictor of the agent's performance under tasks of different difficulty or other problem settings.

In this chapter, I address this limitation by proposing a new mathematical predictive model to estimate the accuracy of artificial agents over tasks of well-defined (quantifiable) complexities. This introduces considerable improvements on the work presented in previous chapters. For instance, the proposed model makes it possible to avert the perpetual

need to simulate an agent over intelligence tests every time we need to predict its performance under a different problem configuration. The model is derived by introducing notions from algorithmic information theory (AIT) into a well-known (psychometric) measurement paradigm, called Item Response Theory (IRT). The model is flexible and can be applied over different types of problems—especially classification tasks. I demonstrate the model over inductive inference problems, which are presented to the agents in the form of intelligence tests. I then discuss how a lower bound on accuracy can be guaranteed with respect to task complexity and the breadth of its solution space using the proposed model. This in turn makes it possible to formulate the relationship between agent selection cost, task difficulty and accuracy as optimisation problems. Further results indicate how to identify the settings over which a group of cooperative agents will be more or less accurate than individual agents or other groups.

5.2 Introduction and Background

Turing’s famous imitation game (Turing, 1950) inspired a range of attempts to measure the intelligence of artificial agents in the twentieth century. More recently, (Hernández-Orallo, 2000) devised a formal (machine) intelligence test consisting of sequence-completion exercises. Two years later, (Bien et al., 2002) used fuzzy integrals to measure intelligence in machines by calculating a Machine Intelligence Quotient (MIQ). Shortly after, (Sanghi and Dowe, 2003) presented a dull computer program that succeeded in passing a variety of IQ tests, raising many questions on the appropriateness of intelligence tests for machine assessment. In 2004, different types of recognition problems were designed such that humans can relatively easily pass compared to machines. These problems, known as CAPTCHAs (Von Ahn et al., 2008), were also proposed as intelligence tests for artificial agents and bots. In the wake of (Legg and Hutter, 2007)’s general definition of *universal intelligence*, many studies (building on the notion of algorithmic information-theory) were put forward to formally quantify the intelligence of individual AI agents (Hernández-Orallo and Dowe, 2010; Insa-Cabrera, Dowe and Hernández-Orallo, 2011; Legg and Veness, 2013; Hernández-Orallo et al., 2012) as well as AI collectives (Insa-Cabrera and Hernández-Orallo, 2013; Chmait, Dowe, Li, Green and Insa-Cabrera, 2016).

Independently, a series of measurement theories have been proposed in psychometrics and have been applied to human intelligence. Perhaps one of the earliest milestones in human intelligence testing was the development of Thurstone’s letter series completion problems (Thurstone, 1938; Thurstone and Thurstone, 1941) and, more recently, Raven’s Progressive Matrices test (Raven and Court, 1998), which recorded strong correlation with Spearman’s general intelligence factor (Spearman, 1904). More general tests consisting of a variety of evaluation tasks (specifically designed to assess cognitive abilities in humans) were also developed, and they came to be known as “Intelligence Quotient” or simply IQ tests. Examples of such tests are the Stanford-Binet test (Roid, 2003) and the Wechsler intelligence scales for adults and children (Wechsler, 2008). For a more detailed historical background on the different approaches used to measure intelligence in the disciplines of AI and psychometrics refer back to Sections 3.4 and 2.3.

Another mainstream achievement in psychometrics was the development of Item Response Theory (IRT) (Lord and Novick, 1968), also referred to as latent trait theory. Item Response Theory is among the most familiar measurement classes used in psychometrics for evaluating traits, or abilities, and producing accurate rankings from test scores, by applying mathematical models to testing data. In the context of IRT, a trait or an ability might be physical or psychological (cognitive and non-cognitive, e.g., might refer to a personality or behavioural characteristic) (De Ayala, 2013). Recently, IRT was successfully adopted to analyse machine learning models by providing an instance-wise analysis of a series of datasets and classifiers (Martinez-Plumed et al., 2016).

In this chapter, I will show how to adapt models from psychometrics and IQ tests, based on notions from algorithmic information-theory, to artificial intelligence in order to estimate the (cognitive) abilities of artificial agents and predict their accuracies.

5.3 Motivations and Major Contributions

Advances in psychometrics are not yet formally applied for predicting the accuracy of AI agents despite their success in evaluating human cognitive abilities. Dowe and Hernández-Orallo (2012) showed that, despite that the AI discipline adheres to the mainstream concept of intelligence (Gottfredson, 1997), general IQ tests might not be appropriate (in their current form) for evaluating machine intelligence. Such tests rather require formal derivations in their test structure to make them suitable for assessing AI.

In fact, even test batteries that might be suitable for practically evaluating AI show some caveats. For instance, such tests measure an average performance (of one or more abilities) of AI agents over a set of tasks or environments but it is ambiguous how the results from these tests can be used to predict the accuracy of an agent over a particular task complexity without actually administering that task to the agent. In addition to many theoretical studies, empirical studies such as (Insa-Cabrera, Dowe, España-Cubillo, Hernández-Lloreda and Hernández-Orallo, 2011; Chmait, Dowe, Li, Green and Insa-Cabrera, 2016) demonstrated that task complexity and the breadth of its solution space are major factors impacting the performance of artificial agents. Hence, quantitatively predicting the accuracy of artificial agents across different task complexities and solution spaces is clearly an important feature that has not been addressed so far.

Furthermore, intelligence test scores can be unreliable since agents, individual or collectives, usually exhibit non-uniformity or disparity between their performances over different problems or task settings. This has strong implications for:

- (i) selecting agents to solve a task, particularly when there is cost (e.g., processing time, or recruitment fee in case of human agents) associated with utilising agents, and
- (ii) understanding the collective accuracy and the relationship between cooperative agents of different (cognitive) abilities.

By merging notions from both psychometrics and (algorithmic) information theory, I developed a hybrid model to quantitatively estimate the accuracy of AI agents over tasks

of measurable complexities. The model is general and flexible to evaluate several types of abilities and agents types over tasks of varying difficulties. I demonstrate its functionality over a class of prediction and inference problems (thus reflecting the agents' inductive-inference abilities) as this class of problems has frequently appeared in the literature of psychometrics (Gold, 1967; Gottfredson, 1997; Feldman, 2003) and artificial intelligence (Blum and Blum, 1975; Dowe and Hajek, 1997a,b, 1998; Li and Vitányi, 2008; Dowe et al., 2011), and is considered as reflecting one of the principal traits of (artificial) intelligence. Using the predictive model, I will show how to identify agents that can guarantee a lower bound on accuracy with respect to task complexity and the breadth of its solution space. This has further implications by allowing one to formulate the relationship between agent selection cost, task difficulty and agent accuracy as optimisation problems. I infer and analyse settings over which a group of (voting) agents can be more or less effective than individual agents, or other groups of agents. I also identify circumstances that can be counterintuitive to the conclusions drawn from intelligence tests.

The next section outlines some important properties and constraints that the model needs to embrace as discussed in (Chmait, Dowe, Li and Green, 2017). I formally define the model and describe how to measure agent abilities in Sections 5.5 and 5.6 respectively. The model is used in Section 5.7 to analyse and predict agent accuracies over different problem settings and complexities, both in individual and collective scenarios (Section 5.8). In the conclusion, I summarise the main outcomes from this chapter and give some directions for future work.

5.4 Desirable Properties for Assessment

Given a subject (agent) to be evaluated over a task/problem:

1. The model must return a *quantitative* measure (on an interval scale) of the estimated subject's accuracy over this task without the need to administer it to the subject.
2. The accuracy of a subject (its probability of success in solving a task) predicted by the model should be proportional to its (relevant cognitive) ability over that task, and inversely proportional to the difficulty of the task.
3. In order to conform to the *limiting* behaviour of real agents, the model should use the asymptotic minimum (probability of correctly selecting a random guess from the sample space) as a lower-bound on accuracy.
4. The model should be applicable over any task of measurable difficulty (which will be revisited in Definition 5).
5. The difficulty measure should be general enough to accommodate a wide range/variety of tasks.
6. The model should be applicable to any agent type or cognitive system (human or artificial).

Earlier information-theoretic studies on (artificial) intelligence (Hernández-Orallo and Dowe, 2010; Insa-Cabrera, Dowe, España-Cubillo, Hernández-Lloreda and Hernández-Orallo, 2011; Chmait, Dowe, Li, Green and Insa-Cabrera, 2016), and inductive-inference (Solomonoff, 1960; Lempel and Ziv, 1976; Papentin, 1983; Evans et al., 2002), discussed two general dimensions of task difficulty. The first dimension relates to Shannon’s entropy (Shannon, 1948) and the uncertainty and breadth of the solution search space, while the second stems from algorithmic information theory, in particular, the Kolmogorov complexity (Kolmogorov, 1965; Li and Vitányi, 2008) of the task. I take into account both dimensions of difficulty in the design of the model.

5.5 A Predictive Model of Agent Accuracy

Inspired by the *two-parameter logistic model* (Birnbaum, 1968) of Item Response Theory (IRT) (Lord and Novick, 1968) and the approach used in (Halmes, 2013) to the measurement of intelligence, I propose a mathematical model (Chmait, Dowe, Li and Green, 2017) for predicting a subject’s expected accuracy on a given task/problem of measurable complexity.

Definition 5 *Let x denote a task/problem of a (theoretical) difficulty \mathcal{D} such that the solution to x belongs to the alphabet (or solution space) $S = \{s_1, s_2, \dots, s_m\}$. I define the accuracy of an agent with ability $\alpha \in \mathbb{R}^+$ over that task to be:*

$$P_{\mathcal{D},\alpha,m} = \frac{1}{m} + \exp^{-\frac{\mathcal{D}}{\alpha}} \cdot \left(1 - \frac{1}{m}\right) \quad (5.1)$$

which corresponds to the probability of that agent guessing the correct solution to x .

The model defined above has the following important properties:

- For a given task of a (hypothetically) negligible difficulty, the probability of solving this task is $\lim_{\mathcal{D} \rightarrow 0} P_{\mathcal{D},\alpha,m} = 1$.
- The probability $P_{\mathcal{D},\alpha,m}$ of a subject with ability $\alpha > 0$ solving a task is (exponentially) proportional to the subject’s ability²¹, and inversely proportional to the difficulty of the task \mathcal{D} , and the breadth of its solution space $m \in \mathbb{N}^+$.
- When task complexity \mathcal{D} is very high relative to α (or when the subject’s ability α is small), the probability of success $P_{\mathcal{D},\alpha,m}$ converges to a random guess equivalent to $1/m$, which is the asymptotic minimum²².

For instance, on a binary test problem (e.g., coin toss problem with $S = \{\text{Heads}, \text{Tails}\}$) with $m = 2$, an agent with ability α has an accuracy $P_{\mathcal{D},\alpha,m} = 0.5 + \exp^{-\frac{\mathcal{D}}{\alpha}}(0.5)$. When the ability α is close to zero, $P_{\mathcal{D},\alpha,m} \cong 0.5$.

²¹In theory, agents might behave worse than random (a very proficient agent trying to err all the time would get a negative ability). Some models, such as the logistic model in IRT, allow for negative abilities.

²²For a more accurate approximation of the true error, one could replace the parameter m used in the model with the **Delta error** function from Probably Approximately Correct (PAC) Learning (Valiant, 1984) with the aim of selecting the hypothesis \mathbf{h} (from a set of hypotheses) generating a low generalisation error. This can also be extended to account for noise.

For many problems, the theoretical task difficulty \mathcal{D} can be derived from the simplest solution (policy) to the task, and therefore can sometimes be linked to the *complexity* of the (description of the) task, or the complexity of the description of its policy. Consequently, the difficulty of the task can be linked to its Kolmogorov complexity (Kolmogorov, 1965; Li and Vitányi, 2008). Since the Kolmogorov complexity²³ is uncomputable, methods like Levin’s *Kt* complexity (Li and Vitányi, 2008; Levin, 1973) or the Lempel-Ziv (compression) algorithm (Lempel and Ziv, 1976) can be used as practical alternatives (to bound it and possibly approximate it). For the rest of this chapter, I will use the Kolmogorov complexity of the task as a derivation of its (theoretical) difficulty. However, it is important to make the distinction between the two, and clarify that the Kolmogorov complexity of a task is not a *synonym* of its difficulty. For example, one should be aware that the task description complexity does not capture the actual difficulty of the task (it serves at most as an upper bound). The difficulty of the task is related to the complexity of its simplest solution(s), rather than to the description of the task itself.

The suggested model returns the probability of a subject solving a given task of a measurable complexity as a function of its (previously measured) ability. The ability could be defined as a vector of weighted atomic sub-abilities s.t. α is a linear combination of $[w_1\alpha_1 + w_2\alpha_2 + \dots + w_t\alpha_t]$. The model in Equation 5.1 is a simple case of the latter where, for some integer $z \leq t$, the ability $\alpha = w_z\alpha_z$ and $\sum_{j=1, j \neq z}^t w_j = 0$ in $[w_1\alpha_1 + w_2\alpha_2 + \dots + w_t\alpha_t]$.

I will use a formal intelligence test from the literature of AI, namely the C-test (Hernández-Orallo, 2000), to measure an agent’s ability α over an important class of tasks used in the evaluation of intelligence. \mathcal{D} and m are input parameters to the model typically being measured by some earlier assessment or derived directly from the problem. I will refer to the model defined in Equation 5.1 as the “IRT model” for brevity, and also use the terms *accuracy* and *performance* alternately (only) as measures of the probability of success at solving a (cognitive) task.

It is important to note that measuring accuracy in this context assumes that there is a selection of discrete *solutions* to choose from. Of course this reflects some important dimensions of intelligence linked to an agent’s ability to perform, for example, classification and inference tasks. Nevertheless, most business solutions would be plans of action that list what needs to be done and how to divide up the work among team members. For example, such plans would take the form of constructed sets/sequences and appropriate task allocations to the group members that lead to specialisation.

5.6 Evaluating Inference Abilities Using the C-test

The C-test (Hernández-Orallo, 2000) is a compression-based intelligence test that measures the ability of a subject doing inductive-inference and finding the best explanation for sequences of various complexities. It is supposed to reflect the *g-factor* (mainly *Gf* or the *fluid* intelligence) of the evaluated subject. The idea is to evaluate a subject over

²³The Kolmogorov complexity is a non-negative integer. Because of *universality* it must be non-zero, so it is a positive integer, and it is usually greater than one. However, some normalised measures have been proposed where the complexity is normalised to a positive real number.

a series of patterns of increasing incomprehensibilities (or complexities) and record its performance. The complexity of a C-test sequence is formally measured using Levin's *Kt* complexity (Li and Vitányi, 2008) as a practical alternative to (and possibly a rough bound on) its Kolmogorov complexity. Example C-test sequences of *Kt* complexities or incomprehensibilities of 9, 12, and 14 (denoted by k_9 , k_{12} and k_{14} respectively) are:

$$\begin{aligned} k_9 &: y, y, x, w, w, v, \dots && \text{answer} \leftarrow u \\ k_{12} &: a, a, z, c, y, e, x, \dots && \text{answer} \leftarrow g \\ k_{14} &: c, a, b, d, b, c, c, e, c, d, \dots && \text{answer} \leftarrow d \end{aligned}$$

5.6.1 Structure of the test

Given an alphabet $\Sigma = \{a, b, c, d, e, f, g, h, i, j, k, l, m, n, o, p, q, r, s, t, u, v, w, x, y, z\}$, and a sequence $seq = (\theta_1, \theta_2, \dots, \theta_m)$ where each $\theta_i \in \Sigma$, the task consists of predicting the next letter $\theta_{m+1} \in \Sigma$ which correctly completes the sequence.

5.6.2 Reward function

Given a C-test consisting of a collection of test sequences $CT = (seq_1, \dots, seq_x)$ with their corresponding answers (solutions) $S = (\theta_{m+1}^1, \dots, \theta_{m+1}^x)$ and corresponding k -complexities $K = (k_1, \dots, k_x)$, the average score \tilde{r} of an agent π with guesses $S' = (\theta_{m+1}'^1, \dots, \theta_{m+1}'^x)$ over CT is:

$$\tilde{r} = \frac{1}{\sum_{z \in K} k_z} \cdot \sum_{z \in K} k_z \times \text{hit}(\theta_{m+1}'^z, \theta_{m+1}^z)$$

where the function

$$\text{hit}(a, b) \leftarrow \begin{cases} 1 & \text{if } a == b \\ 0 & \text{otherwise} \end{cases}$$

and the complexity of the sequence k_z is used as a weight in order to give more importance to more difficult questions.

5.6.3 Inductive inference and choice of C-test

The C-test score will be used to determine the inductive-inference ability α of a subject, further used as a parameter in the model (Equation 5.1). The reasons for selecting the C-test are:

- Firstly, the test by definition measures an (inductive inference related) ability, in this case the ability of finding the best explanation for a given sequence using induction, and completing the sequence by inferring its next character.
- The test is well formulated and is exclusively defined in computational terms. It generates sequences (tasks) within a range of complexities $7 \leq D \leq 15$, using Levin's *Kt* approximation (Hernández-Orallo, 2000) (as a practical alternative to *Kolmogorov* complexity) in a similar fashion to (Hernández-Orallo and Minaya-Collado, 1998, Sec. 5.4).

- Results from the C-test are highly correlated with those from classical psychometric (IQ) tests (Hernández-Orallo, 2000).
- The test sequences are formatted and presented in a quite similar way to psychometric tests. Hence, the test is not anthropocentric and can be applied to machines in the same way it is applied to humans.
- The idea is that there is typically only one best explanation (one exclusive correct answer) for any of the test sequences, making the results uncoincidental and ideally representative of the testee's accuracy.

5.6.4 Measuring abilities

I define in the next few paragraphs (labelled A to F) some simple agent behaviours to be evaluated over the C-test (Chmait, Dowe, Li and Green, 2017). Note that more advanced algorithms for sequence prediction problems exist. But, since the choice of agent behaviours is not particularly relevant to the validity of the IRT model, the selection is restricted to the agent behaviours described in this section for simplicity.

A) Random agent

Given a sequence seq , a *random agent* π^{rand} randomly uniformly selects a letter from Σ and returns it as its answer θ'_{m+1} (Refer to C-test reward function).

B) Mode agent

Given a sequence seq , a *mode agent* π^{mode} looks for the most repeated or frequent letter(s) in seq to predict the next letter. If more than one letter satisfies the criteria, it chooses the left-most one appearing in the sequence.

C) Min-repetition agent

Given a sequence seq , a *min-repetition agent* π^{mr} looks for the least repeated letter in seq to predict the next letter.

D) Min-distance agent

Given $seq = (\theta_1, \theta_2, \dots, \theta_m)$, agent π^{mind} (where the superscript *mind* stands for minimum “*min*” distance “*d*” without spaces) looks for the minimal alphabetical distance (Definition 6) between all consecutive letters of seq and infers the next letter θ'_{m+1} by adding this distance to seq 's last letter θ_m .

Definition 6 *The alphabetical distance $d(\gamma - \beta)$ between two characters β and γ in an alphabet Σ is equal to the difference between their index positions in the totally ordered set (Σ, \leq) in $\text{mod}|\Sigma|$.*

For instance, the distance between any two consecutive letters in the alphabet is 1, and the distance between the first character a and the last one z is equal to $d(z - a) = 26 - 1 = 25$. So, given a C-test sequence $seq = (\theta_1, \theta_2, \dots, \theta_m)$, agent π^{mind} calculates the distance $d^i := d(\theta_{i+1} - \theta_i)$, following Definition 6, between two consecutive elements of seq for all $i \in \{1, \dots, m - 1\}$ returning a pattern (list) of distances $D = (d^1, d^2, \dots, d^{m-1})$. Then, π^{mind} looks for the minimal alphabetical distance $d^{min} \in D$ as follows:

$$d^{min} \leftarrow \arg \min_{d \in D} \text{freq}(d, D)$$

where $\text{freq}(d, D)$ is a function that returns the rate at which d occurs in D .

Agent π^{mind} finally selects $\theta'_{m+1} \in \Sigma$ such that $d(\theta'_{m+1} - \theta_m) = d^{min}$, as the next letter that completes the C-test sequence.

E) Max-distance agent

Agents of this type have an opposite behaviour to that of *min-distance agent*. Given a sequence $seq = (\theta_1, \theta_2, \dots, \theta_m)$, a *max-distance agent* π^{maxd} calculates the distance $d^i := d(\theta_{i+1} - \theta_i)$ between the consecutive elements of seq for all $i \in \{1, \dots, m - 1\}$ returning a pattern (list) of distances $D = (d^1, d^2, \dots, d^{m-1})$. It then looks for the maximal alphabetical distance: $d^{max} \in D \leftarrow \arg \max_{d \in D} \text{freq}(d, D)$ (using the same definition of $\text{freq}(d, D)$ from above).

A Max-distance agent finally chooses $\theta'_{m+1} \in \Sigma$ such that $d(\theta'_{m+1} - \theta_m) = d^{max}$, as the next letter that completes the C-test sequence.

F) Pattern agents

A pattern agent π^{pt} looks for a repeating distance pattern between the elements of seq and attempts to complete this pattern in order to infer θ'_{m+1} .

To implement this behaviour, the problem is first divided into $m - 1$ tasks denoted as $\{t_1, t_2, \dots, t_{m-1}\}$ assigned to agents $\{\pi_1^{pt}, \pi_2^{pt}, \dots, \pi_{m-1}^{pt}\}$ respectively. Agent π_y^{pt} calculates $d(\theta_{i+y} - \theta_i) \forall i \in \{1, \dots, m - y\}$ and generates a list of distances $D_y = (d_y^1, \dots, d_y^k)$ where $k = m - y$, and $d_y^i := d(\theta_{i+y} - \theta_i)$. Then, agent π_y^{pt} searches for the occurrences of the longest possible pattern in D_y and continues D_y by adding d_y^{k+1} following Algorithm 5 defined below. Finally, agent π_y^{pt} makes its guess θ'_{m+1} for the next letter of seq such that: $d(\theta'_{m+1} - \theta_{m+1-y}) = d_y^{k+1}$.

The scores of the above-defined agents over the C-test are used to measure their (inductive inference) ability α and are plotted in Figure 5.1 along with their corresponding accuracies $P_{\mathcal{D}, \alpha, m}$ generated using the IRT model (Equation 5.1). The agents' abilities were calculated as a function of their C-test scores using $\alpha = \tau \tilde{r}$, where α is the ability of agent π with score \tilde{r} , and $\tau \in \mathbb{R}$ is a fitting parameter selected in such a way to:

- (i) ensure that the agent's moderate accuracy, of $0.5(\max P_{\mathcal{D}, \alpha, m} + \min P_{\mathcal{D}, \alpha, m})$, which is equivalent to $0.5(1 + 1/m)$, falls under the area of *discriminative* task complexities $\int_{\mathcal{D}=6}^{\mathcal{D}=16} P_{\mathcal{D}, \alpha, m}$ [following (Hernández-Orallo, 2000)], and

Algorithm 5 Pattern search algorithm

```

1: Input: a set of distances  $D_y = (d_y^1, d_y^2, \dots, d_y^k)$ .
2:   Begin
3:     Extract the unique elements of  $D_y$ .
4:     Store elements in a list  $U_y$  in order of appearance.
5:     Find the starting index for each substring occurrence  $U_y$  in  $D_y$ .
6:     Store index in vector  $v$ .
7:     if  $|v| > 1$  then
8:        $P \leftarrow D_y(v(1) : v(2) - 1)$   $\triangleright v(i)$  is the  $i$ 'th element of  $v$ 
9:     else if  $|v| \leq 1$  &  $|D_y| > 1$  then
10:       $P \leftarrow D_y(|D_y| - 1)$ 
11:    else
12:       $P \leftarrow D_y$ 
13:    end if
14:     $ind \leftarrow |D_y| - |P| \times |v|$ 
15:    if  $ind > 0$  then
16:       $d_y^{k+1} \leftarrow P(ind + 1)$ 
17:    else
18:       $d_y^{k+1} \leftarrow P(1)$ 
19:    end if
20:    return  $d_y^{k+1}$ 
21:   End
22: Output: the next distance  $d_y^{k+1}$  that continues  $D_y$ .

```

(ii) minimise the mean squared error between the IRT model and C-test scores.

The proposed model nicely estimates the agents' average accuracies illustrated in Figure 5.1 despite the large non-uniformity in their individual behaviours and performances.

5.7 Predicting Agent Performances

While results from the C-tests are all alone interesting, we have no means to extrapolate them or predict the agent performances over different sequence complexities and alphabets (or solution spaces) without re-running the test. However, the expected accuracies of an agent can easily be generated from the IRT model over inference tasks of different complexities. An example is illustrated in Figure 5.2 showing the predicted accuracies of agent π^{mind} (refer to the Section 5.6.4) across different hypothetical (Kolmogorov) complexities D and problem solution spaces m .

For any fixed difficulty \mathcal{D} , the IRT model shows that the difference in accuracy measures $P_{\mathcal{D},\alpha,m_1} - P_{\mathcal{D},\alpha,m_2}$ over two solution spaces $m_2 > m_1$ is:

$$\frac{1}{m_1} + \frac{\exp^{-\frac{\mathcal{D}}{\alpha}}}{m_1} - \frac{1}{m_2} - \frac{\exp^{-\frac{\mathcal{D}}{\alpha}}}{m_2} = \frac{(1 + \exp^{-\frac{\mathcal{D}}{\alpha}})(m_2 - m_1)}{m_1 \cdot m_2}$$

meaning that this difference is greater over smaller $m \in \mathbb{N}^+$. This can be observed in Figure 5.2. For consecutive values of m , $P_{\mathcal{D},\alpha,m} - P_{\mathcal{D},\alpha,m+1} = (1 + \exp^{-\frac{\mathcal{D}}{\alpha}})/(m^2 + m)$,

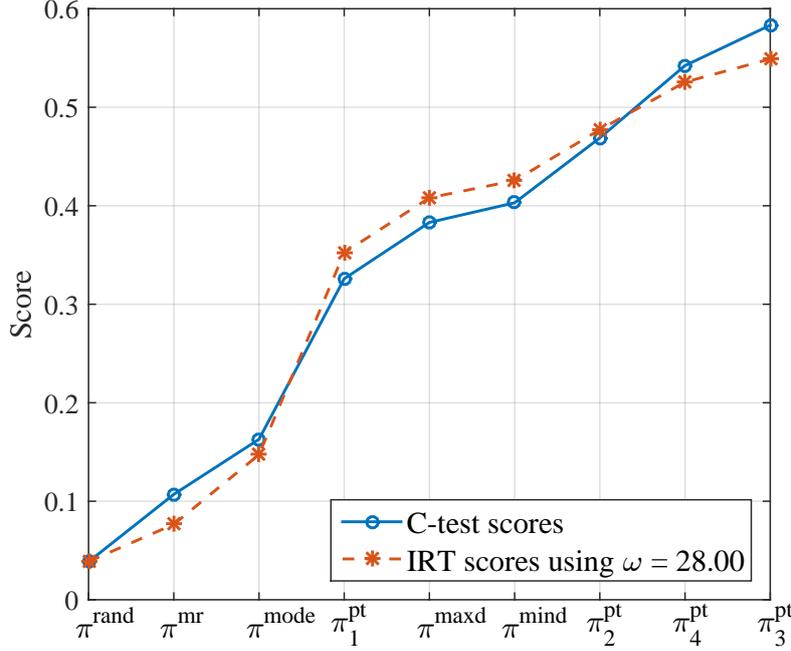


Figure 5.1: Final C-test score \tilde{r} of 9 different agents behaviours (defined above in paragraphs A to F) over a range of task complexities $7 \leq D \leq 15$, and their corresponding IRT accuracies taken from Equation 5.1, using $\alpha = \tau \tilde{r}$ s.t. $\tau = 28$ (Chmait, Dowe, Li and Green, 2017).

and therefore, for very large m , any further increase in m has a negligible effect on the accuracy.

5.7.1 Relationship between accuracy and task difficulty

Figure 5.3 shows the shift in accuracies of a pool of example classifiers of hypothetical (classification) abilities $\alpha \in [1, 8]$ across several \mathcal{D} and m values. We observe that m has a greater influence than \mathcal{D} on the accuracy of those classifiers with poor abilities $\alpha < 3$ and thus their scores are asymptotically bounded by $1/m$, while the opposite is true for more adept classifiers with stronger abilities. This type of analysis can be used to identify the minimal ability value for a classifier to be considered effective compared to, for example, a simple random classifier.

One can further put a bound on the task complexity an agent can solve with a minimal probability of success $P_{\mathcal{D},\alpha,m}$. For instance, if we know m , it is straightforward to calculate \mathcal{D} from Equation 5.1 as:

$$\exp^{\frac{-\mathcal{D}}{\alpha}} = \frac{P_{\mathcal{D},\alpha,m} - \frac{1}{m}}{1 - \frac{1}{m}} \implies \mathcal{D} = -\alpha \ln \left(\frac{m \cdot P_{\mathcal{D},\alpha,m} - 1}{m - 1} \right).$$

Similarly a lower bound on accuracy can be guaranteed with respect to the task complexity and the breadth of its solution space. This is illustrated in Figure 5.4 for agent π^{mind} .

This becomes interesting when a cost function (e.g. processing time, fee) is associated with utilising agents of higher abilities. For example, two agents π_1 and π_2 , with abilities α_1

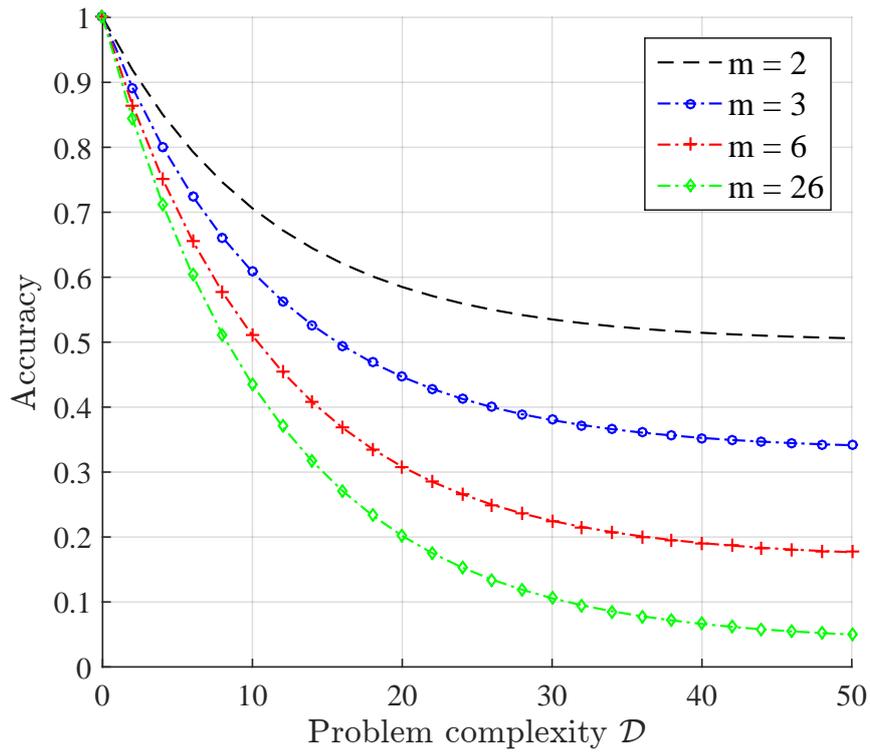


Figure 5.2: IRT accuracy of agent π^{mind} with ability $\alpha = 11.28$ over inference tasks of different hypothetical (Kolmogorov) complexities \mathcal{D} and problem solution spaces m (Chmait, Dowe, Li and Green, 2017).

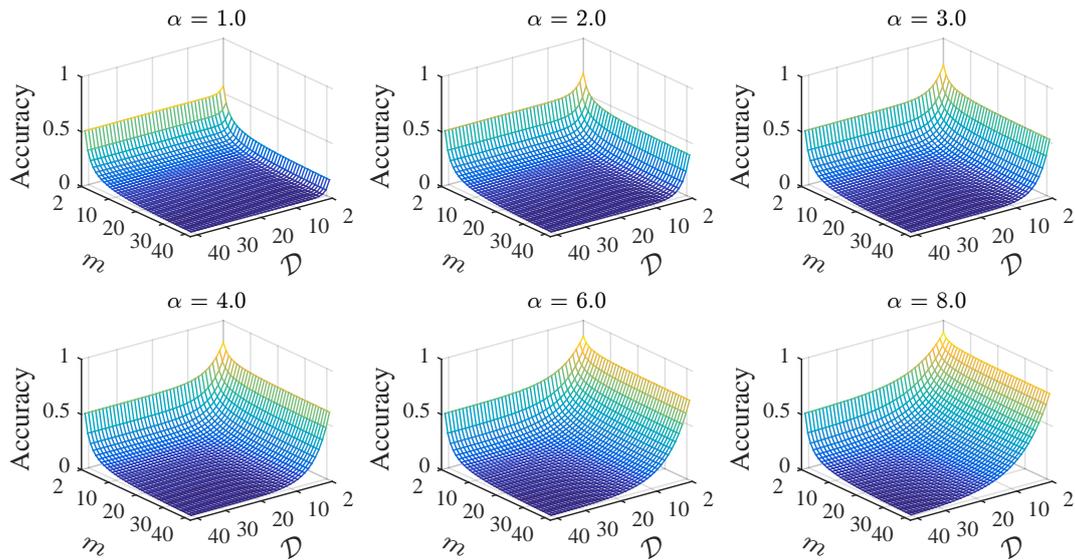


Figure 5.3: Shift in accuracy (from Equation 5.1) across several \mathcal{D} and m values for example classifiers of different hypothetical abilities such that $\alpha \in [1, 8]$ (Chmait, Dowe, Li and Green, 2017).

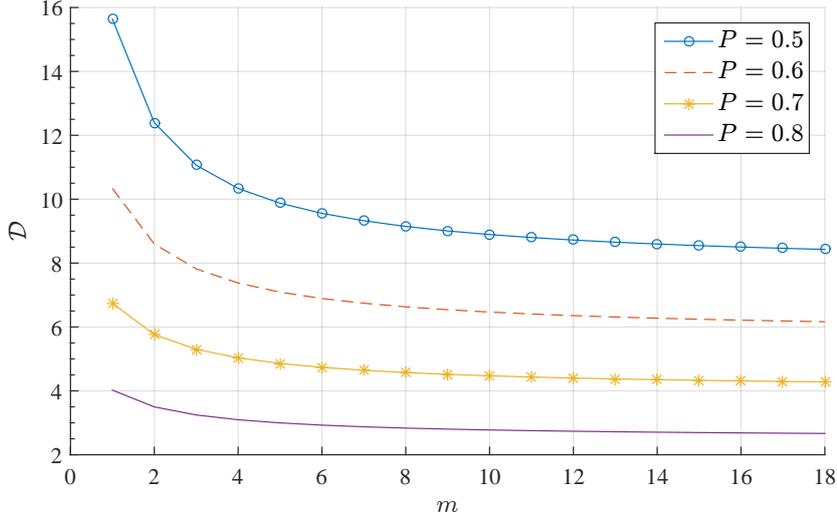


Figure 5.4: Lower bounds on accuracy denoted as $P \in [0.5, 0.8]$ that can be guaranteed with respect to the task complexity D and the breadth of its solution space m for agent π^{mind} with ability $\alpha = 11.28$ (Chmait, Dowe, Li and Green, 2017).

and α_2 and utilisation costs $c_1 = f(\alpha_1)$ and $c_2 = f(\alpha_2)$ respectively, guarantee an accuracy $P_{\mathcal{D}_1, \alpha_1, m} = P_{\mathcal{D}_2, \alpha_2, m}$ under different problem complexities such that $\mathcal{D}_2/\mathcal{D}_1 = \alpha_2/\alpha_1$. If $\alpha_2 > \alpha_1$ (and $c_2 > c_1$) then π_2 can accommodate (a α_2/α_1 factor of) higher problem difficulties with an additional cost $= c_2 - c_1$, while guaranteeing the same accuracy as π_1 . Given a set of tasks of different complexities, a set of n agents of different utilisation costs, selecting the agent to solve these tasks with a *minimum bound on accuracy* of \hat{p} can now be modelled as an optimisation problem: $\arg \min_{1 \leq i \leq n} f(\alpha_i)$, subject to $P_{\mathcal{D}_j, \alpha_i, m} \geq \hat{p}$.

5.7.2 Inferring task difficulty

Alternatively the IRT model can be applied to testing data in order to derive an approximation of the average complexity D of one class of tasks $X = \{x_1, \dots, x_t\}$, assuming the value m for such tasks is already known. For instance, one can empirically evaluate an agent of a known ability α over all task instances $x_i \in X$ and record its average score. Equation 5.1 can subsequently be solved for D using the recorded average score as the accuracy $P_{\mathcal{D}, \alpha, m}$.

5.8 Collective Accuracy of Cooperative Agents

The advantages from adopting the IRT model extend to multiagent scenarios by allowing us to estimate the collective accuracy of a group of agents. For instance, let A be a collective of agents using *simple majority voting* (May, 1952) as a social choice function to elect a solution s_j to a problem x from the set of alternatives $S = \{s_1, s_2, \dots, s_j, \dots, s_m\}$ with only one correct solution $s_i \in S$. Let $Y \subseteq S$: $Y = \{y_1, y_2, \dots, y_n\}$ denote the votes of the agents in $A = \{\pi_1, \pi_2, \dots, \pi_n\}$ respectively regarding their preferred solution to x . When the votes are independent and identically distributed with accuracies p_x , the

probability of collective A finding the solution s_j to x is:

$$P_x(A) = \sum_{k=\lfloor n/2 \rfloor + 1}^n \binom{n}{k} p_x^k (1 - p_x)^{n-k} \quad (5.2)$$

By combining equations 5.1 and 5.2, the probability $P_x(A)$ of a collective of agents $A = \{\pi_1, \pi_2, \dots, \pi_n\}$ electing the correct solution to x with difficulty \mathcal{D} , and alphabet m , using simple majority voting²⁴ becomes (Chmait, Dowe, Li and Green, 2017):

$$P_{\mathcal{D},m}(A) = \sum_{k=\lfloor n/2 \rfloor + 1}^n \binom{n}{k} P_{\mathcal{D},\alpha,m}^k (1 - P_{\mathcal{D},\alpha,m})^{n-k} \quad (5.3)$$

This means that the probability of collective A solving a problem x is the sum of probabilities where at least 50% of the agents are correct. According to Condorcet’s jury theorem (Shapley and Grofman, 1984), $P_{\mathcal{D},m}(A)$ is monotonically increasing when the IRT accuracy $P_{\mathcal{D},\alpha,m} > 0.5$ and vice versa. If A is a group of three agents with unequal accuracies of 0.55, 0.55, and 0.63, its accuracy can be calculated from the agents’ independent choices using majority voting as the probability of at least 2 out of 3 agents finding the correct solution (Kuncheva, 2004, Chapter 4): $(0.55^2 \times 0.37 + 2 \times 0.45 \times 0.55 \times 0.63 + 0.55^2 \times 0.63) = 0.6144$. In fact, even when the agent abilities are unknown, we can still place tight (upper and lower) bounds on the group’s majority voting accuracy as demonstrated in (Matan, 1996) and further simulated in (Halmes, 2013, Sec. 5). The accuracy of an agent collective can thus be sometimes inferred from its agents’ individual accuracies using the IRT model. This allows for measuring the performance of groups of agents and comparing them to individual ones.

5.9 Analysing Individual and Group Accuracies

The accuracy of agent π^{mind} along with the accuracy of different agent collectives (A^1 , A^2 and A^3) over different task complexities and solution spaces are illustrated in Figure 5.5.

We observe that adding agents of equivalent accuracies to the voting process (Collective A^1) improves the accuracy of the group over all tasks where the individual accuracy $P_{\mathcal{D},\alpha,m} > 0.5$ (resulting in the *wisdom of the crowd* phenomenon), while the opposite is true for $P_{\mathcal{D},\alpha,m} < 0.5$.

The key question here is, when is a (voting) collective more efficient than a single agent? To answer this, I calculate the *cut-off point* $\cap_{Y,Z}$ between two evaluated subjects Y and Z . To calculate $\cap_{\pi,A}$ —where the accuracy $P_{\mathcal{D},\alpha,m}$ of an agent π , and $P_{\mathcal{D},m}(A)$ of a collective A , are both equal over some task of complexity \mathcal{D} —we look for the value of \mathcal{D} at which $P_{\mathcal{D},\alpha,m} = \frac{1}{m} + \exp^{-\frac{\mathcal{D}}{\alpha}} (1 - \frac{1}{m}) = P_{\mathcal{D},m}(A)$, which leads to $\mathcal{D} = -\alpha \ln \left((P_{\mathcal{D},m}(A) - \frac{1}{m}) / (1 - \frac{1}{m}) \right)$. If all the agents have similar accuracies (Collective A^1 ,

²⁴More sophisticated voting rules exist such as Borda count, harmonic rule, maximin and Copeland. Such rules require the subject to output a concrete ranking over all possible alternatives of the test/task which inhibits our ability of making exact predictions. Yet, one can still analytically place min and max bounds on team accuracy using different sampling techniques. For state-of-the-art ranking methods over various voting rules and how they compare to the plurality rule refer to (Jiang et al., 2014).

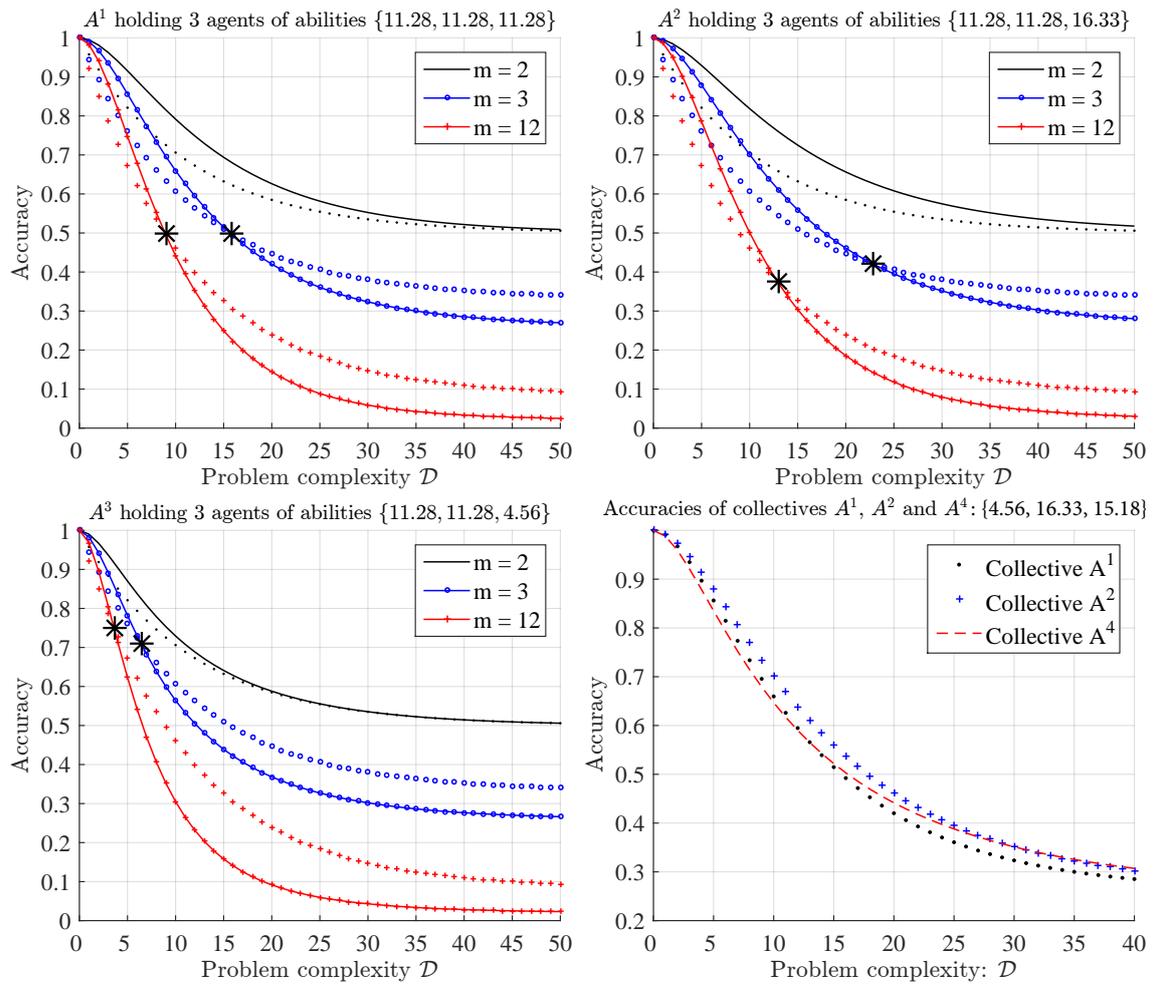


Figure 5.5: Collectives accuracies aggregated using majority voting (Chmait, Dowe, Li and Green, 2017). The accuracy of agent π^{mind} is also depicted as dotted markers in the backgrounds of the first three plots for comparison. The * symbol denotes the cut-off point where the accuracy of π^{mind} meets the corresponding group accuracy.

Figure 5.5), then according to Equation 5.2, they are only equally accurate when $P_{\mathcal{D},\alpha,m} = P_{\mathcal{D},m}(A) = 0.5$ leading to a $\mathcal{D} = -\alpha \ln \left(\left(\frac{1}{2} - \frac{1}{m} \right) / \left(1 - \frac{1}{m} \right) \right) = -\alpha \ln \left((m-2)(2m-2) \right)$. For example, the cut-off point \cap_{π^{mind}, A^1} between π^{mind} with $\alpha = 12.0094$ and A^1 over a problem with $m = 3$ occurs at a $\mathcal{D} = -12 \ln \left(\frac{1}{4} \right) = 16.64$, which can also be verified from the graph.

The cut-off point not only returns the setting over which $P_{\mathcal{D},\alpha,m}$ and $P_{\mathcal{D},m}(A^1)$ are equal, but also illustrates the relationship between the complexity of the problem \mathcal{D} and the size or breadth of its solution space m , with respect to the accuracy of the evaluated group. In other words, the cut-off point indicates the problem complexities and solution spaces over which a collective is more effective than its individual agents. In most real world scenarios voting agents have different abilities and consequently different accuracies. Replacing a group member by another of higher/lower accuracy (Figure 5.5 top-right/bottom-left) improves/diminishes the performance of the group by a measurable amount. For instance, let $A = \{\pi_1, \pi_2, \pi_3\}$ be the group of agents with abilities $\alpha_1, \alpha_2, \alpha_3$ and IRT accuracies (abridged as) p_1, p_2, p_3 respectively over some task x . If the agents' individual votes are independent, the probability $P_{\mathcal{D},m}(A)$ of A correctly guessing the solution to task x by majority voting is:

$$p_1 p_2 (1 - p_3) + (1 - p_1) p_2 p_3 + (1 - p_2) p_1 p_3 + p_1 p_2 p_3. \quad (5.4)$$

When $p_1 = p_2 = p_3$, then Equation 5.4 is equivalent to Equation 5.2. If $A' = \{\pi_1, \pi_2, \pi'_3\}$ is the group of agents with accuracies p_1, p_2, p'_3 respectively s.t. $p'_3 > p_3$, then its accuracy increases by $P_{\mathcal{D},m}(A') - P_{\mathcal{D},m}(A) = p_1 p_2 (p_3 - p'_3) + (1 - p_1) p_2 (p'_3 - p_3) + (1 - p_2) p_1 (p'_3 - p_3) + p_1 p_2 (p'_3 - p_3) = (1 - p_1) p_2 2(p'_3 - p_3)$ since $1/m \leq p_1, p_2 \leq 1$ by definition (Equation 5.1). For $p_1 = p_2 \neq p_3$ the cut-off point $\cap_{\pi_1, A}$ occurs at $\mathcal{D} = -\alpha_3 \ln \left((p_3 - \frac{1}{m}) / (1 - \frac{1}{m}) \right)$ when $p_3 = 0.5$. As a result, we can measure the rise/drop in accuracies of A^2 and A^3 illustrated in Figure 5.5 top-right/bottom-left. For example, for tasks of $m = 3$, \cap_{π^{mind}, A^2} (Figure 5.5 top-right) occurs at a $\mathcal{D} = -16.33 \ln \left((0.5 - \frac{1}{3}) / (1 - \frac{1}{3}) \right) = 22.64$.

5.9.1 Comparing agent collectives

Scores from standard intelligence tests provide us with some sort of scale or ranking of performances of evaluated individuals or groups. Nonetheless, these performance measures might not be valid over certain settings. We observe in Figure 5.5 that voting collective A^1 is more efficient than A^4 (holding agents with abilities $\{4.56, 16.33, 15.18\}$) over inference tasks of $\mathcal{D} < 14$, whereas (counter-intuitively) A^4 scores higher than A^1 over the C-test ($0.51 > 0.38$). However, the opposite is true for tasks of higher complexities. Such scenarios might create confusions as they are frequently encountered and cannot be disclosed from standard intelligence tests.

We also observe that for highly complex tasks with $\mathcal{D} > 25$ collectives A^1 and A^2 record very similar accuracies since $P_{\mathcal{D},m}(A^1) - P_{\mathcal{D},m}(A^2)$ becomes very small. This is coherent with real world observations (although it cannot be drawn from intelligence test

scores) as the accuracy of a subject, or a group of subjects, over extremely hard tasks is likely to converge to a random guess (asymptotic minimum).

5.9.2 A fictitious example

Alice is the CEO of a large organisation. She recently recruited 63 employees to work on a new project. These employees were selected from a large set of applicants after achieving high scores on the job's entrance test, so they all have equivalent levels of abilities and skills required for this type of project. There are several (classification) tasks that need to be solved for this project to be successful.

The 63 employees are given the freedom to decide how they want to get organised and cooperate. Not having worked together before, they decide to work independently on the project-tasks and then go into a majority vote to achieve consensus on the solution for each task. However, there is a disagreement between the employees on how to arrange themselves, e.g., assemble into one group or divide into subgroups of different sizes. The employees know that their accuracy is likely to decrease when working on tasks of higher complexities. Thus, the members decide to simulate their performance over some dummy tasks before beginning their work on the project. They find that in some circumstances they perform very well when voting as one group as a whole, whereas, in others, they score better when they break up into subgroups. After trying different group arrangements and organisation topologies they can't seem to figure out or agree on which one is better.

Bob, one of the group members who works in the field of information theory, decides to assign theoretical difficulty/complexity weights on the set of dummy tasks. He then convinces the other members to repeatedly solve these tasks using three different group arrangements where the team is:

- (i) fully-connected: all members vote as one group,
- (ii) divided into three subgroups of 21 members each. Members of each subgroup vote between one another separately before the three subgroups vote together,
- (iii) arranged in a binary tree-like structure of depth $k = 6$, where each two *children* members vote on the tasks with their *parent* member in a hierarchical bottom-up order until the votes reach the *root* member (final stage).

Bob analyses the resulting accuracy of the group voting in all three types of arrangements (given by the probability mass function of the binomial distribution in a similar fashion to Eq 5.2 used in Sec. 5.8), plots the results (Figure 5.6), and makes some interesting observations. He finds that the group performs well on easy tasks in all three arrangements. However, things become more interesting as task complexity increases. The same selection of agents, using the same decision-making protocol (majority voting) achieves different results over different task complexities. For instance, groups arranged in a tree structure (voting hierarchically) are more effective than others over easy and average task complexities but the opposite is true when the complexity increases. This is not only related to the size of the group, but also to how its members are linked. For example, the fully

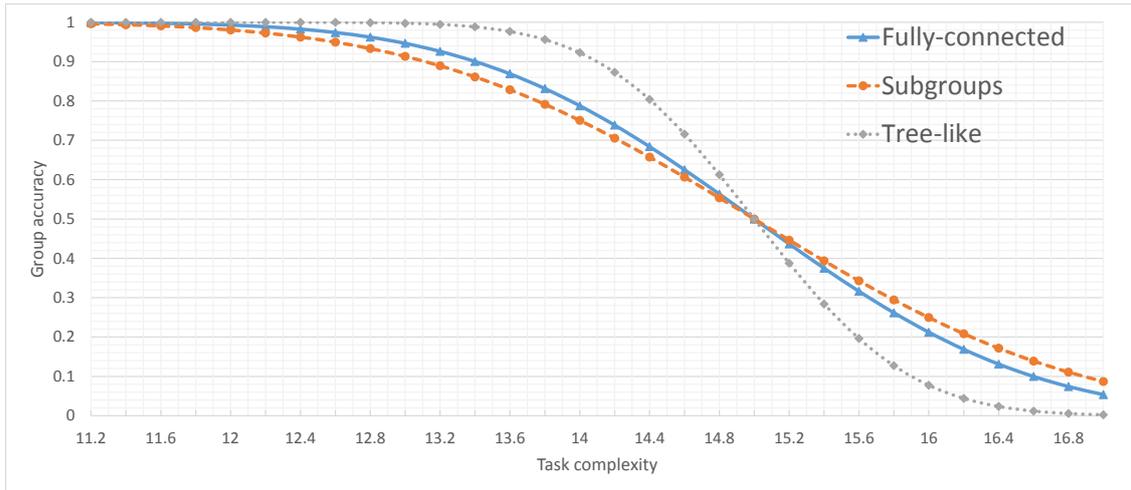


Figure 5.6: Collective accuracy of 63 (majority) voting team members across different group arrangements or topologies. All members are assumed to have equivalent individual probabilities of success over the range of tasks complexities. The individual accuracies decrease (by a hypothetical value of 0.01) as the task complexity increases (by 0.2).

connected group is the largest and yet it does not achieve the highest accuracy (on tasks of hypothetical complexity less than 15). However, a large fully-connected group might perform better when voting on highly difficult tasks. Thus, depending on which approach or criterion for decision-theory is preferred (e.g., Optimist/Maximax, Pessimist/Maximin or the Opportunist/Minimax), one might implement a different group arrangement.

The use of majority voting and the above-described group arrangements are simplistic and used for illustrative purposes. More sophisticated strategies and group topologies can be implemented depending on the context of the problem in hand. As mentioned earlier in this chapter, most business solutions would be plans of action that list what needs to be done and how to divide up the work between team members as for example by performing appropriate task allocations to the group members that lead to specialisation. The intuition here is to show that we can quantitatively analyse decision-making in real business using information-theoretic models such as (extensions of) the one proposed in this chapter.

5.10 Conclusion

Most of the current techniques used for measuring the intelligence of artificial agents return an average performance of these agents over a set of cognitive tasks. While an intelligence score is an accurate measure of an ability associated with the evaluated agent, it is an unreliable predictor of the agent's performance over tasks of well-defined complexities and other problem settings.

In this chapter, I proposed a new mathematical model that is general and flexible enough to predict the accuracy of various types of agents, of different abilities, over different problem settings, therefore fulfilling goal **(G02)** identified in the introduction of this thesis. I illustrated the relationships between the accuracy (and ability) of an agent, the difficulty

of the assessment task and the breadth of its solution space. Using the predictive model, I showed how to identify agents that can guarantee a lower bound on accuracy with respect to task complexity and the breadth of its solution space. Furthermore, I analysed the settings over which a group of agents (using a voting protocol to aggregate their outputs) can be more or less effective than individual agents, or other groups of agents. For instance, I directly inferred from the model the complexity at which a group is expected to record a similar accuracy as an individual agent, and beyond which a single agent is more effective than the group. I also measured the effect (on accuracy) of introducing agents of higher or lower abilities to a group of agents. Finally, I identified possible circumstances that are somehow counterintuitive to the conclusions sometimes drawn from intelligence tests (Sec. 5.9.1). This occurs when a group of agents scores higher than another on an intelligence test yet fails to outperform this same group over some task complexities.

This is a continuing work with many prospects for improvements. I have provided in this chapter initial steps towards a solution for a precarious problem relevant to the measurement and prediction of intelligence. Since the proposed model is properly suited to handle classification problems, my aim is to extend this work in the future to address important topics in this area regarding label noise, classifier dominances, “rough sets” and approximations of decision classes. This can significantly enhance our understanding of questions relevant to task hardness, noise handling, outliers, meta-learning, etc. A more accurate approximation of a subject’s accuracy true error will be sought by using a variant of the Probably Approximately Correct Learning theory (Valiant, 1984).

With respect to group accuracies and objective *Obj06* described in Chapter 1 in particular, more sophisticated voting rules will be used in my future work to analytically reason about team accuracy by analysing the outcomes from different sampling techniques over the agents’ ranked votes. This is likely to have strong implications on a wide range of research disciplines (e.g., health, finance and forecasting services) where machine learning classification plays an important role, as well as other areas of AI such as game-playing and machine intelligence testing.

Chapter 6

Network Science and Intelligence

It may someday happen that the fields of artificial and human intelligence will grow closer together, each learning from the other.

—Douglas K. Detterman, A challenge to Watson (2011)

The research in this chapter will provisionally be published in the following article:

- Nader Chmait, David L. Dowe, David G. Green and Yuan-Fang Li (2017). *Coping with complexity: a multi-disciplinary survey on collective intelligence and its measurement*, Currently under review.

6.1 Overview

As the preceding chapters have shown, it is evident that different kinds of cognitive systems can display some sort of intelligent behaviour at the level of collectives. This has been demonstrated both empirically, using an experimental approach to run multiagent simulations (Chapters 3 and 4), and more formally by mathematical predictive modelling (Chapter 5).

Despite the thorough discussion on intelligence and its main building-blocks given so far, links are still missing to connect fundamental characteristics that are shared between non-uniform types of agent collectives, operating in vastly different environments. This gap raises a series of important questions about the emergence of intelligence. For example, it is not clear whether there is a general model or formalisation that can be used to describe intelligence among the wide range of contexts. Just when, and why, is the whole more than the simple sum of its parts? What are the general principles that allow us to conceive the *emergent behaviour* (Dampier, 2000) in a unified and formal way, across any system or entity type?

In this chapter, I present a new perceptive abstracting tasks and environments using network science. The intention is to facilitate building connections between different research disciplines and cross-fertilising diverse areas of study, ranging from industry and

large organisations to fundamental biology. Hopefully, I can achieve this goal by providing common grounds for all these disciplines to be represented, studied and quantitatively compared with one another.

6.2 What Lies at the Heart of Collective Intelligence?

At first glance, if we are to generalise and identify similarities between the large number of studies surveyed in Chapter 2, we observe that collective intelligence is an abstraction of the complex structure and operation of (possibly large numbers of) simple units—whether they consist of human beings, animals or machines. Aware or not, conscious or not, by committing to (common and sometimes distinct) goals, these simple units demonstrate a high level of differentiation and specialisation, so increasing both their individual and collective efficiency. Here, modularity seems to be a central feature. Whether that is a colony of ants foraging for food, migrating bird flocks, a group of people contributing in online tasks (e.g., Amazon Mechanical Turk) or a bio-inspired artificial multiagent system, by connecting components, interacting, and exchanging resources, intelligent groups emerge from these simple or specialised units.

In addition, we observe common trends in Piaget’s components of Cognitive Theory (Piaget, 1952), explaining how children construct mental models of the world, and the principles underlying the collective behaviour in different cognitive systems. For instance, *adaptation*, a process of learning from schemata²⁵ or modules, relying on *equilibrium*, *assimilation* and *accommodation*, enables the transition in state (or, in other words, the adjustment to the world) in order to deal with new objects or situations, promoting intellectual growth. This interesting concept of adaptation, using equilibrium and accommodation, is pertinent to many of the case studies we have discussed so far. For example, in Section 2.2.1 I mentioned that stability and adaptability are key features in **swarms** (Garnier et al., 2007; Millonas, 1994; Bonabeau et al., 1999). In economy and finance a similar effect is perceived as part of the **efficient market hypothesis**²⁶ (EMH) (Malkiel, 1989) where markets have recipes (or basic building blocks of *intelligent* behaviour in analogy to schemata) to adjust to new information instantaneously, as for example when *adapting* asset and share prices to new information. Moreover, **human crowds** can *adjust* and *accommodate* quickly to changes in their environments by sharing recommendations and aggregating their opinions (e.g., through the web). Therefore, the aggregations of individual human inputs acts as a way of organising knowledge at the collective level allowing them to make accurate predictions and reason about their dynamically changing environment. Adaptability and learning capacity are also key properties of socially intelligent **artificial agents** (Conte, 2002).

²⁵A schema [*plural* schemata or schemas] is defined as “a cohesive, repeatable action sequence possessing component actions that are tightly interconnected and governed by a core meaning” (Piaget, 1952). Piaget also describes schemata as ways of organising knowledge as they constitute the basic building blocks of intelligent behaviour enabling humans to form a mental representation of the world (Piaget, 1952).

²⁶Markets will tend to be quite efficient, but (Dowe and Korb, 1996) points out that the uncomputability of Solomonoff-Kolmogorov complexity puts a limit on market efficiency.

Consequently, we can argue that collective intelligence can be seen as the ripple effect of properties such as *modularity* and *adaptability*, which lead to structure in the central attempt to *cope with complexity*. For instance, modules, just like schemata, are units of knowledge relating to some aspect (component or concept) of the world/environment. These two properties, being deep-rooted in most types of collectives, provide us with a flexible language to talk about their intelligence. Nevertheless, one still needs to resolve the challenge of comparing intelligence between non-uniform types of agents, possibly belonging to different cognitive systems, and operating in vastly distinct environments. In that respect, two directions can be explored:

1. Develop a universal intelligence test which we can administer to any type of entity or cognitive systems on tasks of varying complexities.
2. Seek a flexible measure of complexity that can be used to evaluate each cognitive system or entity in its own environment.

Despite many attempts to design a universal intelligence test (Hernández-Orallo and Dowe, 2010; Hernández-Orallo, 2000; Insa-Cabrera, Dowe, España-Cubillo, Hernández-Lloreda and Hernández-Orallo, 2011; Insa-Cabrera, Hernández-Orallo, Dowe, España and Hernández-Lloreda, 2012; Insa-Cabrera, Benacloch-Ayuso and Hernández-Orallo, 2012; Insa-Cabrera and Hernández-Orallo, 2015), the hypothetical existence of such a test is still highly controversial (Dowe, 2013, Sec. 4.4, p. 23). In fact, such tests must be highly adaptive, and a universal test may simply not be feasible as pointed out in (Dowe and Hernández-Orallo, 2014). Dowe and Hernández-Orallo (Dowe and Hernández-Orallo, 2014, Table 2) surveyed a list of tests designed to evaluate various kinds of subjects, and showed that no test satisfied the universality property. Moreover, we would expect from the wide variety of testing environments and apparatus discussed and analysed in Chapter 2 to make a strong point against the validity and practicality of a universal intelligence test hypothesis.

To clarify this point, I have replotted in Figure 6.1 (without caption details) some of the (artificial) environments appearing in Figures 2.3 to 2.7 of Chapter 2, which simulate (real-world) scenarios and problems from the natural and artificial worlds. These scenarios are mainly an abstraction of searching for a (moving) target (e.g., food or reward) while avoiding injury, and compressing or learning a pattern (e.g., the change in the state of a predator/prey or some resource of interest). If we are to compare performance between agents in ways that can be evaluated in the environments depicted in Figure 6.1 in an unbiased and fair setting, then we must ensure that they are being evaluated over similar environments and tasks. But this task is clearly (almost always) physically not possible. While some natural systems have evolved to adapt and search for solutions over diverse environmental settings, we cannot expect to evaluate and compare ants, birds, AI agents, and human groups in common environments or even on similar tasks. Take for example the tests from Woolley's 2010 experiments (Woolley et al., 2010). The tests consist mainly of 'brainstorming, group matrix and moral reasoning, planning a shopping trip and group typing' activities. Certainly, most of these exercises are strictly associated with human

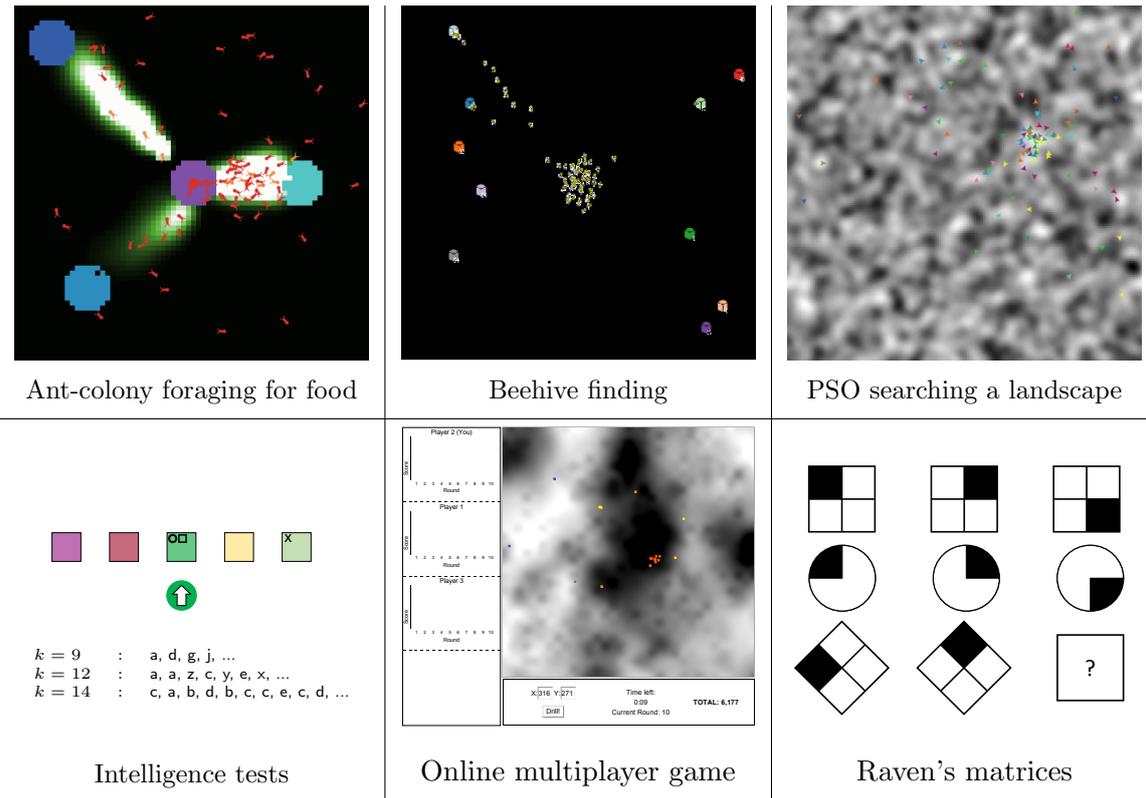


Figure 6.1: Example testing environments whose details were described earlier in Chapter 2, Figures 2.3 to 2.7.

abilities and cannot be administered to other cognitive systems. Also, according to the *No Free Lunch* theorems (Wolpert and Macready, 1997), no system can be expected to perform better than any other over the set of all possible problems. Therefore evaluating entities across a large number of environments, or over *all possible tasks* of a single environment/problem, is somewhat pointless.

Hence, again it seems only plausible to look for a general model to quantify complexity across different environments and tasks, and then evaluate collective intelligence as the ability (or effectiveness) of each kind of system coping with complexity in its own environment. This direction, while not very straightforward to implement, seems to be more reasonable than devising a universal intelligence test. If coping with complexity is a core feature of collective intelligence, then one primary objective would be to define and measure complexity across different environments. A promising path for analysing and understanding (a general measure of) complexity is by resorting to *Information Theory*.

6.3 Information Theory, Complexity and Intelligence

Following our earlier discussion of single agent intelligence in Section 3.4, we note that Bettencourt (Bettencourt, 2009) explored the “formal properties of information aggregation as a general principle for explaining features of social organisation”. This was one of the earliest attempts to quantify the capacity for (collective) intelligence using an information-theoretic approach. It seems that, to produce synergy, it is enough that the

aggregated observations from the members of the collective are all (or sufficiently) statistically mutually independent (Bettencourt, 2009). Bettencourt’s approach is the first step towards quantifying how (and when) information aggregation can sometimes lead to the emergence of intelligence among different individuals, by using the notion of reduction in uncertainty and mutual information (Shannon, 1948). Mutual information also explains some of the mechanism behind Galton’s experiment (Galton, 1907), spatial search and foraging in insects (Dorigo and Stützle, 2009; Bonabeau et al., 1999), the efficiency of markets and opinion polls (Surowiecki, 2005), collaborative filtering and recommendation systems (Goldberg et al., 1992), and so on.

Unsurprisingly, the experimental outcomes from earlier chapters are consistent with Bettencourt’s (Bettencourt, 2009). For instance, Chapter 3 highlights how the principles of uncertainty or entropy (Shannon, 1948) significantly impact the performance of groups of artificial agents. Information aggregation received from interactive agents often leads to groups that can more effectively deal with their environment by using this information (in different ways) to reduce uncertainty and improve their reasoning about the environment. Therefore, **uncertainty** is obviously one dimension of the complexity of an environment. Nevertheless, different experimental results in Chapter 3 reveal that complexity does not exclusively arise from the uncertainty of the environment.

Another major dimension of complexity is how difficult it is **to learn and adapt** to (seemingly random) dynamic environments/tasks. These dimensions of complexity ideally relate to the definition of *intelligence* (Leimeister, 2010; Wechsler and Hardesty, 1964) given in the introduction of this thesis. So what discriminates these two dimensions of complexity? In order to make things clearer, we first need to grasp two fundamental information-theoretic aspects of complexity:

1. the uncertainty or entropy (Shannon, 1948) of a problem, and
2. the algorithmic information-theoretic complexity (AIT) of a problem reflected by its Kolmogorov complexity (Li and Vitányi, 2008)²⁷.

As previously mentioned, on one hand the uncertainty of a problem X with possible states $\{x_1, \dots, x_n\}$ can be calculated as:

$$\mathcal{H}(X) = - \sum_{i=1 \dots n} p_{x_i} \log_2 p_{x_i} \quad (6.1)$$

where p_{x_i} is the probability of event or state x_i of X happening.

On another hand, the algorithmic information-theoretic (AIT) complexity of a string x expressed by its Kolmogorov complexity is defined as the length of the shortest program

²⁷I repeat here a footnote from one of my earlier joint papers. “Note that the concept of Kolmogorov complexity or algorithmic information theory (AIT) is based on independent work of R. J. Solomonoff (Solomonoff, 1960, 1964a,b) and A. N. Kolmogorov (Kolmogorov, 1965) in the first half of the 1960s, shortly followed by related work by G. J. Chaitin (Chaitin, 1966, 1969). For simplicity we use the terms Kolmogorov complexity, or algorithmic complexity, to refer to the algorithmic information-theoretic complexity” (Chmait, Li, Dowe and Green, 2016).

p that outputs x over a reference (prefix-free²⁸) Turing machine U (refer to Definition 1, Sec. 3.6.2).

While the Kolmogorov complexity is uncomputable (Li and Vitányi, 2008), a time-bounded or weighted alternative like Levin’s Kt complexity (Levin, 1973; Li and Vitányi, 2008) or the Lempel-Ziv (compression algorithm) complexity (Lempel and Ziv, 1976; Evans et al., 2002) can be used to bound and possibly approximate it²⁹.

I give a simple example in Table 6.1 illustrating how the uncertainty and Kolmogorov complexity of a problem relate to its difficulty. If we compare the first and last sequences

Sequences	Kolmogorov complexity	uncertainty (in bits)
1, 0, 1, 0, 1, 1, 0, ...	3	$\log_2 2$
1, 0, 1, 1, 1, 0, 1, ...	4	$\log_2 2$
1, 0, 2, 2, 0, 1, 1, ...	6	$\log_2 3$

Table 6.1: Two complexity measures of three example sequence completion problems.

in the table, it is (intuitively) evident that the third sequence is more difficult to complete than the first one. For instance, the third sequence involves picking a solution from a larger set of options or choices. Assume that all possible choices (in the solution space) are equally probable. This means that the uncertainty of the problem X is $\mathcal{H}(X) = \log_2 N$, where N is the total number of choices in the solution space. According to Equation 6.1, the first two sequences have similar uncertainties of $\log_2 2 = 1$ bit (binary options), while the third sequence has an uncertainty of $\log_2 3 = 1.58$ bits. Even when the probabilities are not uniformly distributed, the lowest accuracy one can achieve on the first two sequence problems (after a large number of trials) converges to an asymptotic minimum of $1/2$, while it is $1/3$ for the third sequence.

Besides, if we are to compare the “Kolmogorov” complexities of the first two sequence prediction problems, we find that we need more instructions to describe the pattern generating the second sequence even though both sequences are equally uncertain. I have calculated a bound of the Kolmogorov complexities of these three sequences using Lempel-Ziv’s complexity measure in a similar approach to Section 3.6.2. The resulting complexities are 3, 4 and 6 for sequences one, two and three respectively, showing the importance of this dimension of complexity in describing the difficulty of learning tasks.

²⁸The domain of the Turing machine should be prefix-free. A prefix-free machine is usually considered as being self-delimiting. Such a machine operates with the restriction that the reading head moves only forward and halts when it reads the last bit of the input e.g., it is forced to accept strings without knowing whether there are any more bits written on the input tape (Barrmpalias and Dowe, 2012, Sec. 1).

²⁹I repeat here my comment from footnote 7 of Chapter 3. There is no general algorithm that can determine the Kolmogorov complexity of a given string (Solomonoff, 1964a; Kaspar and Schuster, 1987). Nevertheless, Kaspar and Schuster explicitly recognise in the first paragraph of (Kaspar and Schuster, 1987, Sec. II) that (Lempel and Ziv, 1976) provide an appropriate alternative measure of the Kolmogorov complexity of a string by calculating a number $c(n)$, instead of the length of the program which generates a given string of length n , which is a useful measure of this length. Furthermore, (Evans et al., 2002) discuss why the Lempel Ziv 78 Universal compression algorithm (Ziv and Lempel, 1978) is a computationally efficient method towards approaching the *estimation* of the Kolmogorov complexity.

6.4 Intelligence as Coping With Complexity

Why is complexity important, and how does it relate to intelligence? If we can quantify complexity, then we will be able to measure intelligence in terms of the complexity of the tasks that an individual or a collective can efficiently solve, or more generally in terms of the complexity of the environments in which they can operate. So it is important to distinguish what is considered as a complex environment or task, and what is not. For instance, estimating the number of jelly beans in a jar containing only two jelly beans (very low uncertainty) is pointless. Similarly, the winner from an election in which always the same candidate has been elected for 20 consecutive years is going to be predictable (low algorithmic information-theoretic complexity).

The good news is that we can measure the AIT complexities of numerous classes of problems. These classes are shared across various types of cognitive systems and environments, including for example: inductive inference, compression, pattern recognition and learning, as well as exploration/exploitation and search problems. The key point is that these classes of problems are tightly related to (the measurement of) intelligence (Schaie, 1979; Ryabko and Reznikova, 2009; Chaitin, 2002; Wallace, 2005; Solomonoff, 1964a, 1986; Dowe and Hajek, 1997a,b, 1998; Hernández-Orallo and Minaya-Collado, 1998; Hernández-Orallo and Dowe, 2010; Dowe et al., 2011) and have been used to evaluate diverse kinds of entities (Ryabko and Reznikova, 2009; Reznikova, 2007; Hernández-Orallo, 2015; Hernández-Orallo et al., 2016; Schaie, 1979; Hernández-Orallo, 2000; Hernández-Orallo and Dowe, 2010; Insa-Cabrera, Dowe, España-Cubillo, Hernández-Lloreda and Hernández-Orallo, 2011; Chmait, Dowe, Li, Green and Insa-Cabrera, 2015).

To wrap up, complexity partially arises from the underlying uncertainty of the environment and also the algorithmic information-theoretic difficulty of the task to be performed, which are two well known measures of complexity pertinent to the field of information theory. This bifold dimension of complexity can be used as a tool to analyse and understand systems operating in complex environments and estimate their intelligence. For brevity, I will refer to these two dimensions of complexity by the term *dual criterion complexity* for the rest of this chapter. One of the main advantages of using such an approach is that it captures both the *individual complexity* of a task instance x , e.g. sequence or pattern complexity (using $K(x)$), as well as ensembles and a diversity of those, e.g. $X = \{x_1, \dots, x_n\}$ (using $\mathcal{H}(X)$). These concepts regarding the nature of complexity and its connection to homogeneity, symmetry and disorder are very important and have been previously discussed in biology (Kauffman, 1993), mathematics (Waldrop, 1993), physics and chemistry (Prigogine and Nicolis, 1989), computer science (Li and Vitányi, 2008), etc.. For more details on the intuition behind using a dual criterion complexity measure and its implications I also refer the readers to (Papert, 1983) in which the author discusses the notion of complexity in the context of information theory and formal language theory.

6.4.1 Measuring and comparing intelligence across various cognitive systems

Following the above formulations, one way to calculate the uncertainty in an environment is by looking at the number of all possible states in which resources/rewards might be allocated (neighbour solutions which can be reached) in that environment. The complexity of a task might consist of the AIT complexity of the canonical or cognitive activities to be performed³⁰. For example, this might correspond to the complexity of the movement pattern of the predator/prey to escape/trace, the complexity of the path leading to resources like food, or the seeming randomness in a sequence of numbers in an IQ test. The 2D Λ^* environment used in Chapter 3 is one of many sample environments whose complexity can successfully be measured following the dual criterion complexity. Nonetheless, approximating these complexity measures is not always straightforward and can be problematic in the case of large and dynamic environments. It requires a thorough study of the environment and the task to be performed.

To show the importance of using a general model to measure the difficulty of a problem (based on the notions of uncertainty and algorithmic complexity), I refer back to the test interfaces illustrated in Figures 2.4, 2.5, 2.6, and 2.7 of Chapter 2. All of these interfaces are significantly different from one another. Yet, all of them could be administered to human beings while guaranteeing an equivalent measure of difficulty. For instance, Levin's *Kt* complexity (Levin, 1973) (a computable version of Kolmogorov complexity) has been used by (Insa-Cabrera, Dowe, España-Cubillo, Hernández-Lloreda and Hernández-Orallo, 2011) and (Hernández-Orallo, 2000) to measure the complexity of the two tests illustrated in Figures 2.4 and 2.5 respectively, although these tests have totally different interfaces and presentations. Other studies such as (Dowe and Hajek, 1997a, Sec. 2 and 4)(Dowe and Hajek, 1997b, Sec. 2)(Dowe and Hajek, 1998, Sec. 2 and 4) and (Strannegård et al., 2013; Ragni et al., 2011) also show how (approximations of) the Kolmogorov complexity could be applied to many other IQ tests, including Raven's Progressive Matrices depicted in Figure 2.6.

Despite significant differences between test structures and their dynamics, it is interesting that information theory allows us to evaluate insect colonies, artificial multiagent systems as well as human groups across different tasks of the same difficulty. Assuming that we have a clear understanding of the evaluation tasks, we can begin by identifying a series of steps that can be used to evaluate and compare the (collective) intelligence of (different) cognitive systems. A basic description of these steps is captured in Figure 6.2 and will be further developed throughout the rest of this chapter. Consequently, we can

³⁰One might wonder why not use classes from computational complexity theory like \mathbb{P} and \mathbb{NP} to measure the problem complexity instead. Intuitively speaking, different problems/tasks belonging to the same complexity class (or even more, different instances of the same problem) can have different complexities in terms of uncertainty and Kolmogorov complexity (e.g., depending on the problem instance size). Secondly, while in nature there exists collectives like swarms with the capacity of solving intractable problems in real-time, it is infeasible (due to, for example, time constraints) to evaluate individuals or even small groups over such intractable problems.

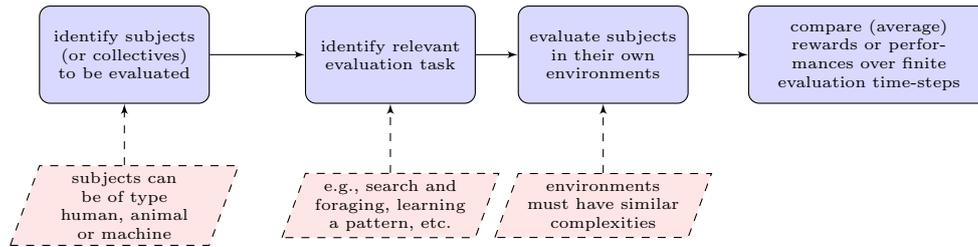


Figure 6.2: A high-level diagram showing a series of steps that can be used to evaluate and compare the intelligence of entities of different cognitive systems. The diagram will be developed incrementally in Figures 6.3 and 6.6.

try to answer questions related to the ones raised earlier in Section 6.1 like “how to quantitatively analyse and compare the intelligence of different kinds of agents operating in distinct environments?”.

6.5 Network Theory and Intelligence

At this point, it is clear that the complexity of the intelligence test inhibits the agents’ (group) performances. However, it is also desirable to link the theory of complexity to some real world problems, and show how the results from the science of complexity might be used to cross-fertilise different research areas.

It is known that numerous real world problems (e.g., path-finding, scheduling, constraint satisfaction problems, ...) and environments (e.g., social networks, computer networks, molecular structures, maps, motifs and patterns, ...) can be efficiently represented in terms of network or graph problems. In fact it has been shown (Green, 1994, 2000; Green and Bransden, 2006; Green, 2011) that a *network structure* is inherent in every model used to represent complex systems, and also inherent in the state space of every automaton or array of automata. Consequently, this indicates that a large number of natural and artificial environments are likely to have an underlying network structure.

Just as we can measure the complexity of strings and sequences, so can we measure the complexity of networks and graphs. Many studies have applied information-theoretic complexity (Kolmogorov complexity and also Shannon’s entropy) to measure the complexity of graphs (or networks) such as (Bonchev, 1995, 2004; Mowshowitz and Dehmer, 2012; Hearn, 2006) and (Li and Vitányi, 2008, Chapter 6).

In order to compare the performance of collectives of different types or cognitive systems, one can first model the environments in which these collectives operate as network or graph models, and calculate their complexities following an information-theoretic approach similar to the dual criterion complexity measure. Each type of collective (human, animal or artificial) could then be repeatedly evaluated in its own environment, and its rewards averaged and compared to those of other collectives or individual agents evaluated in environments (networks) of similar complexities. These ideas are introduced into the diagram illustrated in Figure 6.3 as a refinement of the one given earlier in Figure 6.2. In the next section I give examples of common network structures and show how network science can help us evaluate and understand the behaviour of various cognitive systems.

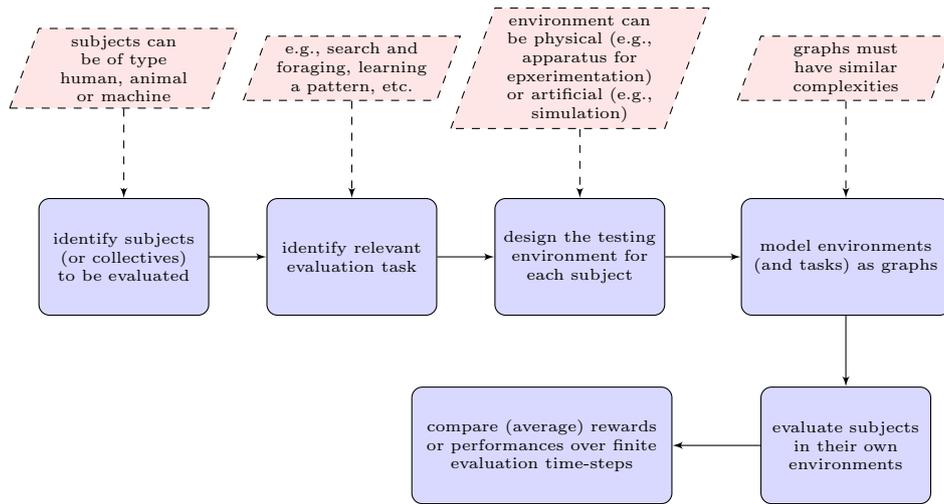


Figure 6.3: A refinement of the diagram from Figure 6.2 introducing a couple of new intermediary steps related to the design of the testing environment and its representation as a graph. A more elaborate diagram is given in Figure 6.6

6.6 Integrating Network Science and Real World Disciplines

Several topologies are found to be common in networks like the ones illustrated in Figure 6.4. These topologies might be associated with social networks like small-world net-

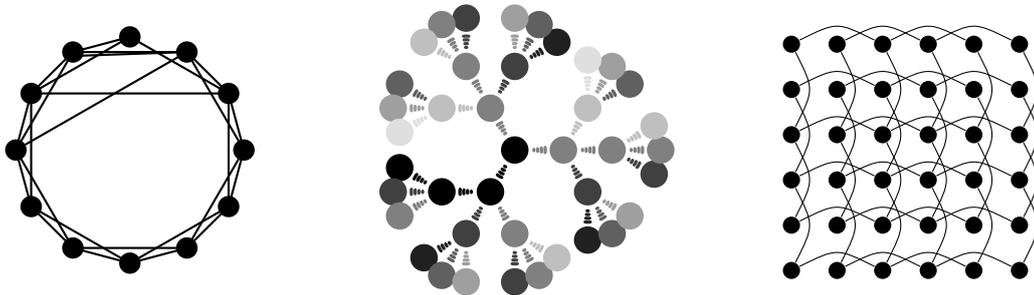


Figure 6.4: Different types of network structures or topologies representing to small-world, tree-like and toroidal environment spaces (from left to right).

works (Watts and Strogatz, 1998) commonly found on Twitter (left-hand side of Figure 6.4), hierarchical and modular networks often represented by trees or fractals (centre of Figure 6.4), scale-free networks (Barabási and Albert, 1999), as well as other natural and artificial environment structures like mazes and toroidal environments (right-hand side of Figure 6.4).

Following the above abstractions to the environment and task complexities, we can quantitatively measure and compare the effectiveness of insect colonies foraging for food to that of a group of cooperative online users, or a group of artificial agents searching a fitness landscape for rewards, and therefore cross boundaries between human, animal and machine entities. This is only possible because we have the opportunity to model real world environments (and tasks) as network structures, and measure their complexities.

Definition 7 Modelling environments as graphs. Given an environment μ and a set of environment states $Q = \{q_1, q_2, \dots, q_n\}$ of μ , the function $\pi(q_i, q_j) = 1$ if μ can subsequently be in the two states q_i and $q_j \in Q$, and q_i and q_j are called neighbour states. The set of all pairs of neighbour states of μ is $N \subseteq \{(q_i, q_j) \mid q_i, q_j \neq i \in Q \text{ and } \pi(q_i, q_j) = 1\}$. Environment μ can be modelled as a graph $G = (V, E)$ where $V = Q$ and $E \subseteq \{(q_i, q_j) \mid q_i, q_j \in V \text{ and } (q_i, q_j) \in N\}$.

For instance an environment can be encoded as a graph $G = (V, E)$ with vertices $V = \{1, 2, \dots, n\}$ and each vertex represents a (different) state of that environment. This graph can be further encoded as a string $S(G)$ of length $n(n - 1)/2$, where each character corresponds lexicographically to an edge connecting two vertices such that $S(G) = e_{1,2}e_{1,3} \dots e_{1,n}e_{2,3}e_{2,4} \dots e_{n-1,n}$, where in turn the edge connecting vertices x and y $e_{x,y} = 1$ if $(x, y) \in E$ or zero otherwise³¹. An example of this encoding scheme is illus-

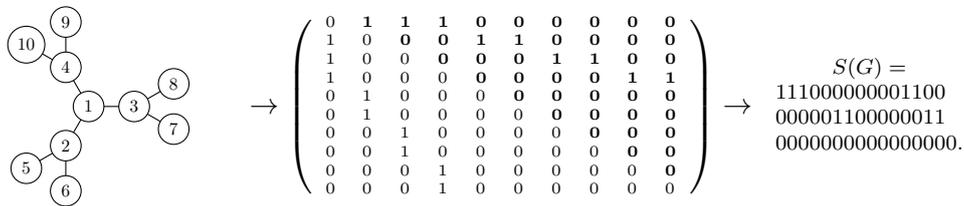


Figure 6.5: Encoding graphs as binary strings over two steps. The resulting string is used for estimating the (Kolmogorov) complexity of the graph.

trated in Figure 6.5. Now imagine a dynamic task over this example environment where agents must forage for food or warm shelter such that a reward r_i (food, convenient temperature, etc.) is dynamically allocated to a (possibly) different environment state (graph vertices) at time-step t_i . For example assume that at $t_1, t_2, t_3, \dots, t_{11}$ rewards can only be received when at states “1, 4, 1, 3, 1, 4, 9, 4, 9, 4, 1” of Figure 6.5 respectively. This can easily be generalised to larger time-steps as the same (circular) rewarding pattern can again be allocated for the following eight time-steps and so on. For example if the test runs for 25 time-steps, the rewarding (task) pattern P can also be represented as the string $S(P) = \text{“1413149494114131494941141”}$, whose Kolmogorov complexity could in turn be approximated.

I next give an example of how to calculate the Lempel-Ziv complexity (Lempel and Ziv, 1976) of a string which is considered as an appropriate measure of the Kolmogorov complexity. I follow an approach similar to (Kaspar and Schuster, 1987, Sec. II) and

³¹Note that this string encoding scheme is dependent on the labelling order of the vertices. Thus, isomorphic graphs might have different binary representations and consequently different Kolmogorov complexities. However, a bound on complexity (and randomness deficiency) can be found for different (labelled) graphs as discussed in (Li and Vitányi, 2008, Chapter 6) and (Hearn, 2006, Chapter 2). We label vertices following a breadth-first technique in order to exploit modularity and capture the pattern underlying the (structure of the) environment which enhances the graph’s compressibility. In fact, most patterns in nature (tree leaves, rivers, etc.) are incrementally generated from a source (e.g. root of a tree) analogously to a breadth-first cross-walk. Breadth first search techniques have also been shown effective for mining frequent patterns (Vaarandi, 2004), pruning and maximal pattern mining (Aggarwal and Han, 2014, Chapter 2) and other graph pattern matching problems (Eppstein, 1995). Also note that, while encoding graphs as strings is possible, there are alternative computational models that are foundationally based on graphs rather than strings.

(Lempel and Ziv, 1976, Sec. II). The general idea is to check whether the newly inserted part of a given string is contained in the vocabulary of (or can be reconstructed by copying) its previous parts (substrings). Consider a string $S = s_1 s_2 \dots s_n$. I will construct S starting from an empty string E . I denote by s_r a newly inserted digit to E that is not obtained by simple copying (of previous parts of E). At each of such insertions I check whether “ s_r ”, “ $s_r s_{r+1}$ ” and so on up to “ $s_r \dots s_n$ ” can be constructed by copying one of E 's substrings. For example, let $S = 001100$,

1. I start by inserting the first digit $E = [0;]$ so the last non-copied digit $s_r = s_1 = 0$. I append a semicolon to E for ever new s_r .
2. $s_1 \dots s_r = 0$ (from above $s_r = s_1 = 0$), $s_{r+1} = 0$, $s_{r+1} \in (s_1 \dots s_{r+1-1})$, so $E = [0; 0]$
3. $s_1 \dots s_r = 0$, $s_{r+1} s_{r+2} = 01$, $s_{r+1} s_{r+2} \notin (s_1 \dots s_{r+2-1})$, so $E = [0; 01;]$, and $s_r = s_3$ (append semicolon)
4. $s_1 \dots s_r = 001$, $s_{r+1} = 1$, $s_{r+1} \in (s_1 \dots s_{r+1-1})$, so $E = [0; 01; 1]$
5. $s_1 \dots s_r = 001$, $s_{r+1} s_{r+2} = 10$, $s_{r+1} s_{r+2} \notin (s_1 \dots s_{r+2-1})$, so $E = [0; 01; 10;]$, and $s_r = s_5$ (append semicolon)
6. $s_1 \dots s_r = 0011$, $s_{r+1} = 0$, $s_{r+1} \in (s_1 \dots s_{r+1-1})$, so $E = [0; 01; 10; 0]$

The Lempel-Ziv complexity is equal to the number of substrings between the semicolons (the use of bracket in “ $E = []$ ” is only to improve readability). In this example the complexity of $S = 001100$ is equal to 4. Similarly, we can measure both environment and task complexity for the previous example illustrated in Figure 6.5. For instance the environment μ modelled as a graph G has a complexity of $K(G) = 7$ while the task corresponding to learning the pattern P has a complexity $K(P) = 6$.

At this stage, we have a formal method enabling us to model many types of problems and environments that are relevant to various cognitive systems, measure their complexities, and use them to evaluate and compare the (collective) intelligence of these systems. This method is illustrated in the refined diagram appearing in Figure 6.6.

6.7 Comparing Graph Complexities

Now that we have a method to model environments and tasks a graphs and measure their complexities, I will briefly discuss the complexity measures of some example graph topologies abundant in nature. Figure 6.7 shows six different graphs, G_1, G_2, G_3, G_4, G_5 and G_6 , along with their measures of uncertainty $\mathcal{H}(G_i)$ and Kolmogorov complexity $K(G_i)$, in addition to a normalised measure of $K(G_i)$ denoted by $\hat{K}(G_i)$ used as an indicator of the formation of a pattern in the graph (revisited later).

We observe from Figure 6.7 that graphs containing a larger number of vertices $|V|$ record a higher uncertainty measure $\mathcal{H}(G_i) = \log_2 |V|$. For instance, it is easier to search an environment with a small number of states than a large one. Moreover, graphs generated using the same mechanism (or Turing machine, like G_1, G_2 , and G_3) record a similar

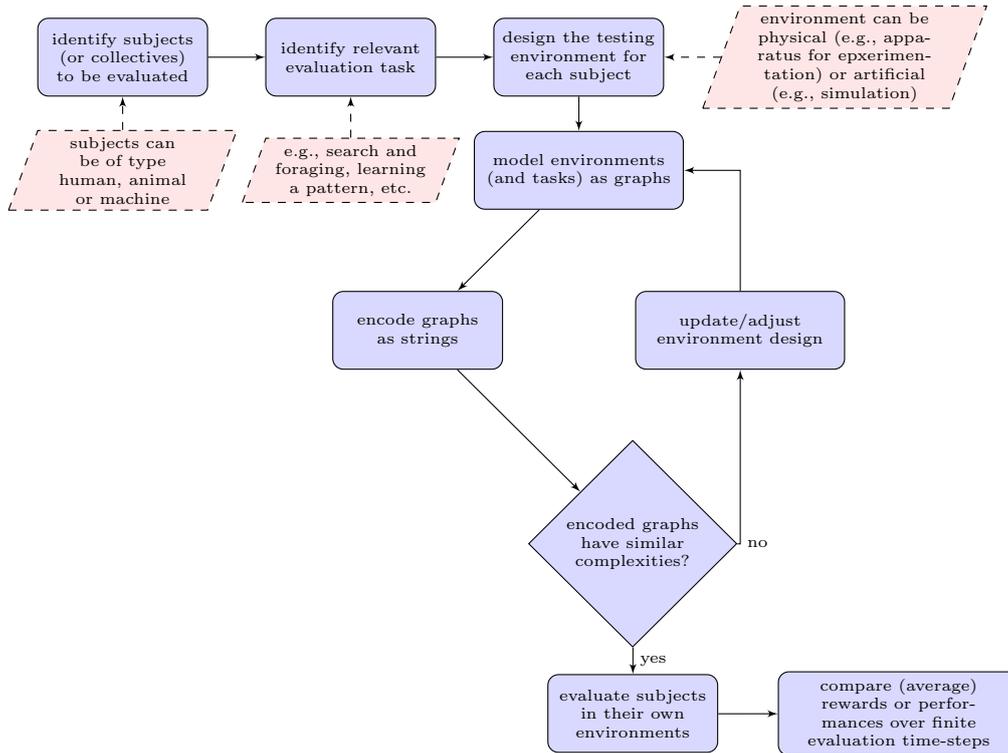


Figure 6.6: A diagram summarising a preliminary methodology for to evaluating and comparing the (collective) intelligence between (different types of) agents operating in vastly distinct environments. Environments are modelled as graphs and their complexities are measured. Agents are then assessed over equally complex environments.

Kolmogorov complexity despite the difference in their number of vertices. On the other hand, the more compressible the graph is, the lower its $K(G_i)$. We can intuitively see that it is easier (or it requires less information) to describe the complete graphs G_1, G_2 , and G_3 than describing G_5 . Similarly, it is much easier to describe the scale-free graph G_5 and the helix (spiral) graph G_6 than describing a random graph such as G_4 . This is a remarkably useful technique to compare environments and graphs, but there is still one additional attribute that is not expressed in our measures: *modularity*. That is to say how much pattern is replicated during the formation of these graphs.

Lempel and Ziv (Lempel and Ziv, 1976) showed that the asymptotic behaviour of their complexity measure $LZ(n)$ for a random string of length n converges to a value $\lim_{n \rightarrow \infty} LZ(n) = \tau(n) \equiv n/\log_2 n$. Moreover, for a $h = -[p \log_2 p + (1 - p) \log_2 (1 - p)] \leq 1$ where p is the source entropy that can have a maximum value of 0.5 (e.g. probability of finding a 1 in a random binary string), they show that $\lim_{n \rightarrow \infty} LZ(n) = h\tau(n)$. This means that a deviation from the value of $\lim_{n \rightarrow \infty} [LZ(n)/\tau(n)]$ might capture the formation of a pattern in the string (Kaspar and Schuster, 1987), e.g. in the case where $\lim_{n \rightarrow \infty} [LZ(n)/\tau(n)] < h$.

Therefore, I will also introduce a normalised measure of complexity $\hat{K}(G)$ to hopefully capture the appearance of a pattern in a graph G . This is calculated as follows: given a graph G encoded as a binary string $S(G)$ of length n , the normalised complexity $\hat{K}(G) = K(G)/(n/\log_2 n)$. We can observe from Figure 6.7 that the more frequent

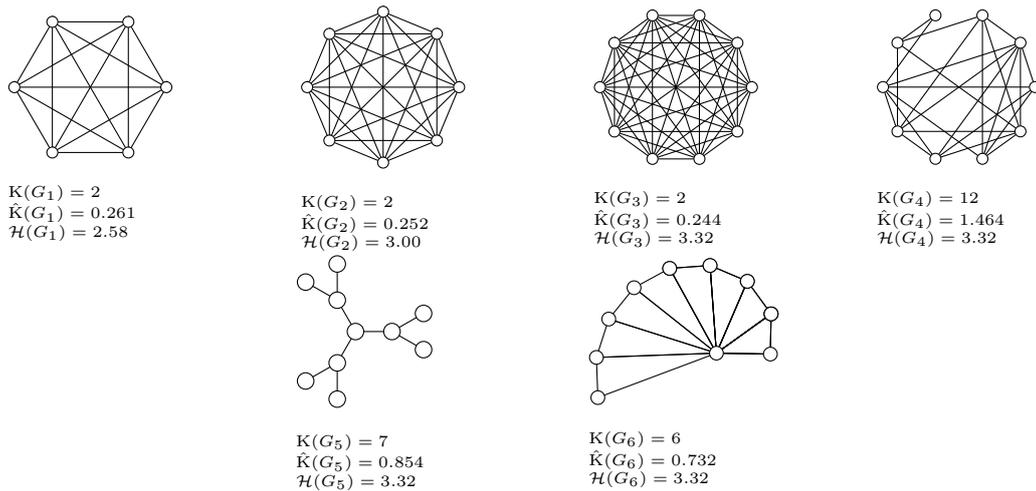


Figure 6.7: Different network structures and their uncertainties $\mathcal{H}(G_i)$, as well as an approximation of their original and normalised Kolmogorov complexities $K(G_i)$ and $\hat{K}(G_i)$ respectively, measured in terms of their Lemple-Ziv (normalised) complexities.

a pattern generating the graph is, the smaller is its normalised measure of complexity. For instance, $\hat{K}(G_1) > \hat{K}(G_2) > \hat{K}(G_3)$. Modularity, as discussed before, is an important characteristic of environments especially with regards to the emergence of intelligence in these environments.

Finally, it is worth noting that the proposed complexity measures that are applied to different network structures (e.g., those in Figure 6.7) ignore the connection costs in the network. As we will see later in Section 6.9, some theories suggest that natural environments evolve following a structure that minimises their connection costs. For example, the formation of a hierarchical tree structure in a network of n vertices is more frequent than a mesh structure since (in a full binary tree) it only requires $n-1$ connections, whereas $n(n-1)/2$ connections are needed in a fully connected network.

Moreover, in order to capture the complexity of elaborate, more sophisticated, real world environments, one might consider introducing new details and dimensions to the above-described measures of complexity. This can still be achieved using information theory. For example, another dimension of complexity related to the topology of a network and its spatial organisational structure can also be captured from graphs. This is illustrated in Figure 6.8 which depicts the topological complexities of all simple connected graphs of five vertices, where each graph represents a chemical molecule. The complexity of each graph is a combination of both its Kolmogorov complexity and its Shannon entropy, with additional weights ascribed on the (entropy) probabilities of each node (atom or subelement) of the graph (molecule) (Bonchev, 1995). These probabilities reflect the distribution of atoms composing the molecule into classes according to their chemical nature. Similarly, such graphs could alternatively be adopted to model and measure the complexity of other types of environments consisting of five spatially distributed states.

With regards to maze-like environments and labyrinths, (Hillier, 2007, p. 91-98) has investigated how to measure the difficulty of a maze or, in other words, the *maze-iness* of an environment. This was done by developing a measure (in the range $[0, 1]$) called

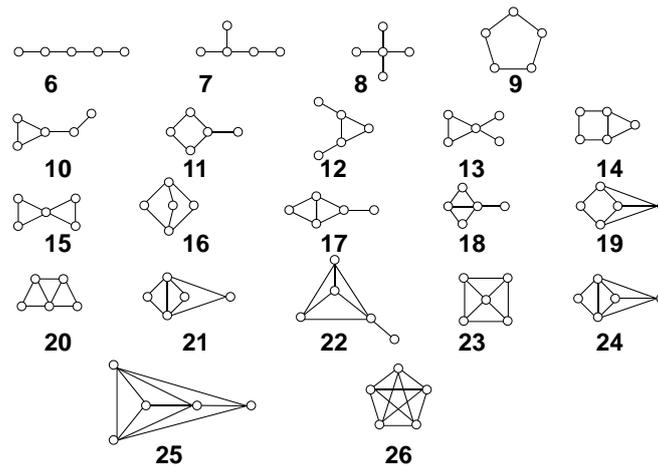


Figure 6.8: All simple connected graphs of five vertices and their *topological* complexities calculated following their Kolmogorov complexity and Shannon’s entropy using a weighted probability distribution (referred to as magnitude-distribution) on the elements making the graph [Redrawn after (Bonchev, 1995)].

intelligibility, which basically detects the “relationship between what is immediately visible from a single location in a maze/housing estate/neighbourhood and how accessible that same place is from other locations in the area” (Dalton and Dalton, 2017). This measure of *intelligibility* could thus be applied to network graph in the same way it can be applied to real geographical areas or neighbourhoods.

Other information-theoretic studies based on the notion of Minimum Message Length (MML)³² (Wallace, 2005) have discussed formal approaches to the inference of networks and the measurement of complexity of network graphs, such as (Comley and Dowe, 2003, 2005), (Wallace, 2005, Sec. 7.4) and (Ooi and Dowe, 2005; Dowe and Zaidi, 2010; Visser et al., 2012), which might also help capture additional dimensions of complexity—as might also carrying out uncertain reasoning/inference in probabilistic logic networks (Goertzel et al., 2008).

6.8 Sample Test Problems and Their Encoding

In Section 6.6 I discussed a preliminary methodology that can be applied to model real world search and learning problems as graphs, and gave examples in Figures 6.5 and 6.7 of how this can be achieved. In this section I give examples of common and important real world problems, and show how they can be modelled as graphs and encoded using the proposed scheme.

6.8.1 The scheduling problem

Scheduling problems appear very frequently in a wide range of our everyday activities. A classic scheduling problem is exam timetabling. The main idea behind timetabling is to

³²See footnote 9

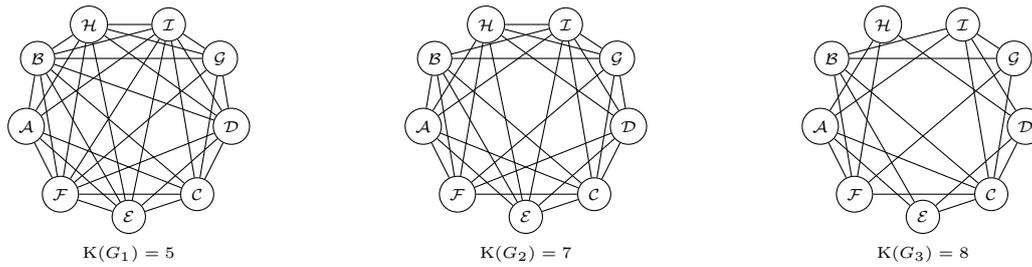


Figure 6.9: Graph models G_1, G_2 and G_3 of the three timetabling examples 1, 2 and 3 respectively, appearing in Figure 6.9 and their Kolmogorov complexities $K(G_i)$.

avoid conflicts when scheduling the exams (e.g., ensuring that exams taken by a common student are allocated at different times).

Table 6.2: Three instances of exam timetabling problems of different number of conflicts.

	Example 1 (4 conflicts)			Example 2 (9 conflicts)			Example 3 (14 conflicts)											
	A	B	C	D	E	F	G	H	I	A	B	C	D	E	F	G	H	I
A				X			X					X	X			X		
B											X		X					X
C							X					X						X
D	X								X	X					X	X		
E							X								X	X		X
F									X		X		X	X				X
G	X			X						X	X		X	X				X
H			X						X	X		X	X			X	X	X
I											X			X	X			X

One possibility is for a timetable to be modelled as a graph with one vertex for each state (exam) and an edge connecting neighbour states with no conflicts. In other words, two neighbouring states (exams) can be scheduled together in the same time-slot. In fact, similar results will be obtained if we choose to connect vertices with conflicts instead (by flipping all ones to zeros in the adjacency matrix) which indicates that this method can be flexible with respect to different interpretations of the nature of the problem or environment state. It is obvious that a table with no conflicts, modelled as a fully connected graph, is the least complex. To highlight how the proposed method captures the essence of the problem difficulty (in terms of how hard it is to schedule due to the existence of conflicts between exams) I give three timetabling examples using the same number of exams to be scheduled (number of vertices) in each example, and only changing the number of conflicts between them. Let $\{A, B, C, D, E, F, G, H, I\}$ be the set of exams to be scheduled. For every two distinct conflicting exams (that must be taken by common students) in Table 6.2 an “X” is added to their corresponding table cell of intersection. Three timetabling examples are given in Table 6.2 using 4, 9 and 14 conflicts for the same set of exams. These examples are further modelled as graphs in Figure 6.9 and their complexities are calculated in a similar manner to Figure 6.5, encoded following the alphabetical labelling order of exams. We observe from Figure 6.9 that the measured complexity of the problem increases with the number of conflicts in the corresponding example.

However, it is important to note that if we keep increasing the number of conflicts in the schedule (given a fixed number of variables/states), the measured complexity of its graph must start to decrease at some point. It is controversial what this implies in terms

of the difficulty of the problem. At first glance, one can interpret this as a limitation of the proposed methodology. Otherwise, this can be explained as a genuine decrease in the problem complexity. For instance, not much information is required to describe a schedule in which almost all exams have conflicts, and consequently the problem might be easily identified as either infeasible or unsolvable. As a matter of fact, this is very similar to the phase transitions associated with the hardness of NP-hard problem instances (Cheeseman et al., 1991). In SAT (propositional satisfiability) problems, a phase transition occurs as the ratio of the number of clauses to variables of a SAT instance changes, showing an *easy, to hard, and again to easy* pattern as this ratio increases (Mitchell et al., 1992; Gent and Walsh, 1994). As the number of constraints increases, this easy→hard→easy pattern also occurs when applying the graph encoding methodology proposed earlier.

6.8.2 The travelling salesman problem: comparing biological and artificial systems

Taking into account all the considerations discussed in this chapter, can we compare bees to artificial agents? Bumblebees (or *Bombus terrestris*) forage for food and visit a vast number of flowers in order to collect nectar and pollen before returning to their nest. While foraging in flowers, bumblebees seek the optimal route (Lihoreau et al., 2012). This is similar to the travelling salesman problem (TSP) (Bellmore and Nemhauser, 1968) where a subject has to find the optimal route covering all cities that minimises his/her trip costs. An example of the TSP is illustrated in Figure 6.10 below³³. In this section, I will show

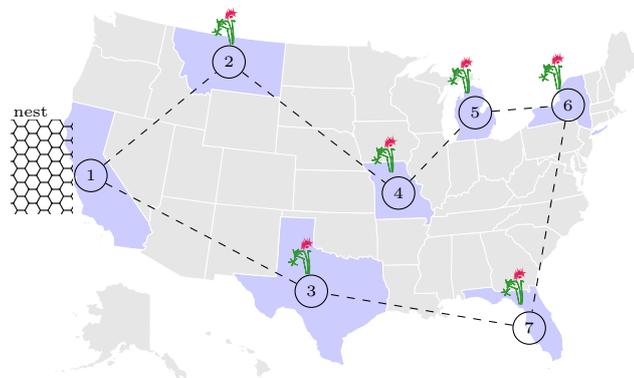


Figure 6.10: An example of the Travelling Salesman Problem (TSP) where a subject, originally located at the state of California, needs to find the least cost route visiting all of states of Montana, Texas, Missouri, Michigan, New York and Florida, and back to his/her original destination, without passing through the same state more than once. The problem is modelled as a graph with vertices denoting the destinations to be visited and edges correspond to a path between one destination and another. A tour (covering all vertices/cities/flowers) is shown in the figure starting from, and back to, California/nest (vertex with label 1).

how to follow the steps in the diagram of Figure 6.10 in order to compare the performance

³³The TikZ scripts used for drawing the USA map and flowers appearing in Figure 6.10 are modified after (D’Ossualdo, 2016) and (Botoeva, 2012) respectively.

of bumblebees to another cognitive system: the ant colony optimisation (ACO) search heuristic.

1. The first step is to identify the subjects to be evaluated. In this example these are the ant colony optimisation (ACO) heuristic (which has been successfully used to provide adequate solutions to the TSP (Dorigo and Gambardella, 1997)) and a collective of bumblebees.
2. The second step is to identify the assessment task. For the ACO heuristic, the task is an instance of the travelling salesman problem. Whereas for bees, originally located at a nest, it consists of finding efficient routes among flowers and back to the nest without pollinating the same flower more than once.
3. The next step is to design the testing environment for each of the subjects. There are many ways in which the TSP can be presented to artificial algorithms which I will not address here. As for bees, one possibility is to examine them over artificial flowers (Lihoreau et al., 2012) placed in a closed field or space.
4. We then encode the testing environments as graphs. For artificial agents, each city is modelled as a vertex in a graph. Cities (vertices) have routes (edges) connecting them. Similarly for bumblebees, a flower is modelled as a vertex in a graph, and flowers (vertices) have a path (edges) between them. A weight can be assigned to each edge connecting two nodes in both graphs corresponding to the distance between the connected nodes. Other weights could also be taken into consideration (e.g., elevation, weather conditions and so on). Depending on the test objectives (e.g., if time is taken as a factor), the ACO heuristic can be fine-tuned such that each ant (artificial agent) requires the same period of time, to travel from one end of an edge of distance d to another, as a bee. As for the assessment task, it can be encoded as the string corresponding to the elements in the set of all the permutations of the vertices making a tour (hence covering all possible solutions).
5. We encode the graphs modelling each testing environment as strings (Recall Section 6.6). This encoding returns a higher uncertainty measure when the number of vertices is larger (cities/flowers). Moreover, the encoding of the task might return a higher Kolmogorov complexity when the problem instance has more edges connecting the vertices (as opposed to having limited paths between the vertices).
6. We measure the complexity of the encoded graphs (environments and tasks) following the methodology described in Section 6.6.
7. If the graphs have the same complexities, we can proceed with evaluation. Otherwise, the environments should be updated (e.g., flowers and/or paths connecting them added or removed) until they are equally complex.
8. The subjects are evaluated over a finite period of time in their own environments.

9. Their performances are finally measured and compared. It is straightforward to record and measure the performance of artificial agents. One way to do that for bees is to have them wear a radar transponder and track them as in (Lihoreau et al., 2012). The finale route corresponding to the ordered set of visited destinations/flowers can be considered as the solution from each of the cognitive systems, and consequently the two systems can be compared with one another while operating in their own environments.

These types of experiments illustrate a new approach to the measurement and comparison of intelligence across different types of entities as their ability of coping with, and operating in, their own (complex) environments. The structure of the evaluated entities and their inter-dynamics are abstracted. While this abstraction can be very useful, it is not straightforward. It rather requires expert and careful study and design before abstracting tasks and modelling environments as graphs³⁴.

6.9 Agents as Networks

In a large part of this thesis we have been looking at collective intelligence as the ability to operate in, and cope with, complex environments. Indeed environments play a fundamental part in shaping our actions and behaviours, and in many cases determine our capacity (as a group/society) for intelligence. Nevertheless, the agents' individual behaviours in turn contribute to the characteristics of the environment they belong to. Many of such examples have been discussed earlier in Section 2.2 in the contexts of human, animal and artificial intelligence. In other words, complex (social, group) structures and environments stem from the interactions, and more importantly the interdependence, between individuals.

In this section I look at collective intelligence as the transition from a micro (agent level) to a macro (social level) intelligent behaviour, and how this helps to cope with complexity. To analyse the emergence of intelligence at the aggregate level from simple individual interactions one can also represent and study interactive individuals as networks.

Here again, the same principles discussed in Section 6.2 seem to apply. For instance, the adaptability of agents to one another, and the way they cluster into modules, are main characteristics of their complexity and intelligent emergent behaviour. A large number of studies have looked at the interdependencies between agents as networks in order to understand their collective behaviour. In social sciences (Skyrms and Pemantle, 2009; Eguíluz et al., 2005; Zimmermann et al., 2001), business and economics (Gatti et al., 2005, 2011) and markets (Weisbuch et al., 2000), patterns have been detected which were captured by small-world (Watts and Strogatz, 1998) and scale-free networks (Barabási and Albert, 1999). A major contribution in the area of collective intelligence was the game-theoretic approach to the formation of (intelligent) networks (Jackson and Wolinsky,

³⁴Alternative types of problems can be modelled as graphs and used for evaluation such as, the stable-marriage problem (Gusfield and Irving, 1989), map-colouring (Barnette, 1983) problems, or other practical (map or Geographical Information Systems) search problems like (Amaneddine and Chmait, 2009).

2003; Colman, 2013, 2014). The analysis of the strategic behaviour of individual agents has explained many of their (intelligent) collective behaviour and decision-making (Oakley, 2010; Rong et al., 2010). More elaborate versions of the prisoner's dilemma, such as the N -person prisoner's dilemma (Colman, 2014, Section 8), have been thoroughly investigated. Namatame and Chen (Namatame and Chen, 2016) recently collected and analysed a broad number of studies on the integration of agent-based modelling and network science. The authors discuss a series of important questions (Namatame and Chen, 2016, Chapter 1) on agent network dynamics that closely relate to the notion of collective intelligence, such as:

- why do agents choose to cooperate, or to be kind, even under the presence of high temptation to defect?
- how do they interact and adapt to each other in such a way to produce aggregate outcomes of interest?
- how can one identify the individual behaviour that contributes to consistency in the behaviour of interest at the macroscopic (collective) level?
- can purely individual actions bring change to the society?

Modelling agents as networks and simulating their behaviour might not necessarily be enough to quantitatively measure and compare their (collective) performances, especially when these agents fall under different types of cognitive systems. In fact, many agent models are considered to only reflect the dynamics associated a particular network (Gil and Zanette, 2006). However, networks dynamics are very accurate indicators of the agents' collective behaviour and have been used successfully to model phenomena like opinion formation, rationality and the public good. For example, opinion and sentiment formation can be represented as simple network-based decision rules (Namatame and Chen, 2016, Section 2.2.4) (e.g., cellular automata) which lead to structure and stability in spite the chaotic initial distribution of opinions.

Just as environments evolve according to some underlying structure, interactive agents themselves have an underlying structure which can be captured via network modelling. This was expressed by Newman (Newman, 2006) in his proposed measure of network modularity. Newman points out that most networks of scientific interest (including those representing social and biological systems) divide naturally into communities or modules and shows how to detect and characterise these community structures (Newman, 2006). Attributes like modularity, complexity and emergence are brought to light as part of our real-world management practices used in group decision making, which in turn have been abstracted as networks. The time-honoured strategies of “divide and rule” and compartmentalisation are frequently used by large organisations in order to deal with complex problems by dividing them into smaller, simpler parts, with separate subgroups allocated to deal with each part (Green, 2014, Chapter 8). At first glance these ubiquitous strategies that shape our social hierarchies and intelligence seem ideal. Nevertheless, there are circumstances in which these well-known strategies fail (Green, 2014, Chapter 8) as a

result of promoting limited closed-box thinking (the assumption that our local everyday activity is a closed box independent of other activities) which is incoherent with the reality.

In other studies the distinction between environments, or otherwise agents, as networks becomes blurry and almost disappears. For example, Mengistu et al. (Mengistu et al., 2016) studied networks at a higher level of abstraction by investigating the evolution of network structures. Results show that “networks without a connection cost do not evolve to be hierarchical, even when the underlying task has a hierarchical structure”, whereas, when connection costs are introduced, networks evolve to be both modular and hierarchical. In addition, modular networks exhibit a higher overall performance and adapt much faster to new environmental conditions. The take-home message from (Mengistu et al., 2016) is that indirect selection is a fundamental factor to the evolution of an organisational structure, in particular hierarchy and modularity, in the main goal of reducing the net cost of network connections. This is an elegant example of how modularity and adaptability promote coping with complexity through evolution by reducing the cost of the network structure.

6.10 From Intelligence Tests to Other Spheres

The study of agents (and their environments) as networks seems to be a promising area that can make a revolutionary leap in understanding group dynamics and their intelligence by connecting the studies of intelligence to other spheres, notably business decision-making. I have discussed in Section 3.12 of this thesis how the effectiveness of agents adopting different organisational and network structures can vary from one structure to another over intelligence tests. A natural extension of this work is to conduct similar sorts of intelligence measurement studies, modelling interactive agents as networks, and further using the outcomes to understand the effect of different types of group organisations, team dynamics, management and business decision-making strategies. Many research problems can therefore be connected to one another, such as those addressing interaction structures of groups (Watts and Strogatz, 1998; March, 1991; Watts, 2004; Easley and Kleinberg, 2010; Mason and Watts, 2012; Anicich et al., 2015) and (collective) decision-making techniques used between the group members (Pentland, 2006; Yu et al., 2010; Krause et al., 2011; Charness and Sutter, 2012)—which can all be evaluated and compared quantitatively. It is indeed an elegant example of how collective intelligence can lead AI to different areas. In other words, we can analyse the characteristics of a particular group strategy across different disciplines and problems, and make conclusions that extend beyond a narrow application or area of research. On one hand, intelligence studies can be used to determine whether observable differences in scores after evaluation are the result of the network topology or otherwise the agent (or problem) type or other underlying dynamics. On the other hand, one can conduct studies of intelligence to explore the relationship between the evolution of a network topology, structure or group strategy and the difficulty/complexity of the task. For instance, given a certain group of agents, we can quantitatively analyse and determine which problem complexities necessitate the use of group strategies such as *divide-and-rule* and *compartmentalisation* that are commonly used

by large organisations (Daft, 2012) (revisited below). In this context, intelligence studies are no more used to only measure intelligence, but rather as a medium to explore effective group dynamics and characteristics, as well as to understand how the effectiveness of business decision-making is influenced by the organisational structure of the agents across different areas of research.

Connecting intelligence studies to other domains is particularly advantageous in *management*. For instance the choice of arrangement to be used between the components of a system, and the management of relationships between them, can be adequately studied. The importance of such choices and their potential impact on group performance were raised by (Malone and Bernstein, 2015) in the context of the *Star Model* organisational design (Galbraith, 2002). The authors outline different categories of decisions about design that managers and system designers need to make in order to enhance business productivity and organisation success. These decisions relate mainly to the group *strategy* and its *structure* (overall goals and objectives, and decision-making technique), as well as its *processes* and *rewards* (flow of information among components and their incentives).

Other examples in (Malone and Bernstein, 2015) show the importance of the *integration* of (outputs from) group components and their *interdependencies* (e.g., in hierarchies and markets) and the relevance of the group's functional and divisional structures (e.g., departmentalisation and geographical divisions of the components). Such organisational strategies, patterns and structures are not confined to management and business practices but also appear in social and political sciences (e.g., scale-freeness and modularity), engineering (e.g. manufacturing modular components), computer science and robotics (e.g., agent communication, cooperation and feedback), bioinformatics (e.g., representation of living organisms and spread of diseases), ecology (e.g., motifs), etc..

Therefore, using the methodology presented in this chapter and a similar experimental approach to the one in Chapter 3, we can compare the intelligence of different types of groups implementing various interaction strategies, (business) organisational structures and management practices by evaluating their effectiveness in dealing with complex problems. I illustrate in Figure 6.11 some of these important (cross-disciplinary) strategies and structures that we can now quantitatively study and analyse in a wider scope than before. The networks in Figure 6.11 represent agents (nodes) interacting with each other according to some arrangement or organisation (defined by the edges/connections between the nodes/agents), sometimes with additional permissions or restrictions assigned to their interactions (e.g., arrows or special notations). The *intelligence* of (possibly the same set of) agents or individuals can be tested by implementing such group arrangements and organisations over relevant problems and environments of equivalent complexities. Consequently, their performances can be analysed and compared in order to understand how each of these arrangements shaped their effectiveness over different types of problems and contexts. For instance, in the context of business organisations and management, sample business operations can be presented in the form of intelligence tests and administered to these agents. The traditional model of decision making in large organisations relies on hierarchies that break down a problems into smaller, specialised parts (e.g., finance,

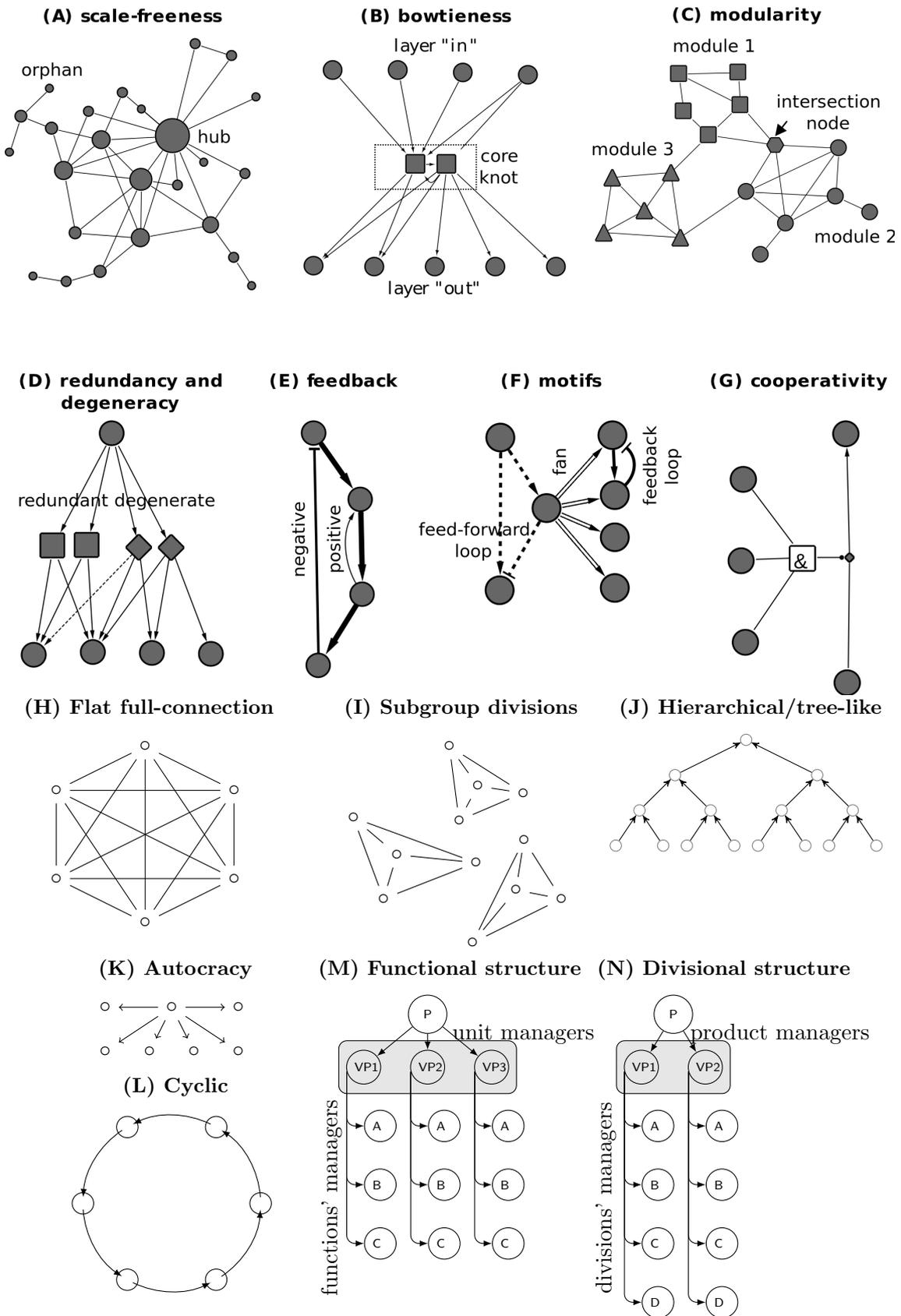


Figure 6.11: Important group structures and templates responsible for robust network properties and distinctive group behaviours. [Subgraphs **A** to **G** taken from (Barillot et al., 2012, Fig. 9.1), an open access figure distributed under the Creative Commons Attribution-Share Alike 3.0 Unported license, which permits sharing and remixing provided the original work is properly cited and attributed to the original authors].

marketing, production), with information fed up the hierarchy and orders flowing down. This is demonstrably the traditional model in business, military and government. These sorts of arrangements are illustrated in subfigures J, M and N of Figure 6.11. For instance, departmentalisation within an organisation allows upper management to monitor and control the activities of its employees efficiently (Daft, 2012). Nevertheless, different organisational structures form depending on whether the business departmentalises according to the activities performed by individual groups within the business (Figure 6.11, subgraph M), or otherwise according to geographical areas, products, services or markets (Figure 6.11, subgraph N), where in this case each division operates as an autonomous (smaller) business (Daft, 2012). Such organisational structures can be assessed by applying intelligence studies on agent groups after restricting the flow of information between the group members in order to model the implemented business structure.

Furthermore, many modern companies augment fixed modules (e.g., departments) with ad-hoc teams composed of experts from different areas, to address particular tasks. However, their plans and/or findings almost always pass as recommendations to higher management or board of directors. Such scenarios can be captured using different network topologies (or connections) between the node/agents, and further by using directed graphs (or other additional constraints on the graph) to model the flow of information and hierarchical levels of decision-making. Finally, the resulting networks (reflecting the group structures of the agents) can be compared over intelligence tests in order to analyse their performances and understand their behaviours across different types of complex problems.

In biology and ecology, agents can be simulated over various artificial environments (of measurable complexities) using different sorts of group structures to understand their evolution and how well they adapt to such environments. In robotics, multi-robot systems can be assessed quantitatively over tasks requiring cooperation (e.g., collective pursuit, mapping a landscape or team games) using different group strategies to compare and enhance their performance.

The above considerations are possible given that this thesis provides common grounds for the evaluation of various cognitive systems operating over different environments and problems, whose complexities can be formally measured. This can open the doors for a new understanding of how the wide range of studies in collective intelligence relate to one another, and provide the opportunity for their quantitative analysis in different contexts.

6.11 Summary and Future Work

Most living entities adapt to survive and flourish. They divide naturally into modular interactive components that demonstrate remarkable intelligence at the large scale. This intelligent behaviour occurs in human societies, insects and other swarms. Consequently, adaptation and modularity resulted in diversified types of interactions taking place between entities and their environments which, in turn, has promoted most organisms to operate in vastly different environments. On one hand, analogies among these organisms and other systems can easily be observed. On the other hand, big difficulties are encountered with regard to making quantitative comparisons between these systems' performances

and their collective behaviours. This inhibits our ability of designing and understanding different kinds of collective intelligence and the concepts shared among them.

This chapter proposed a preliminary approach towards addressing these difficulties. Abstractions and simplifications were made to establish common grounds for comparison. This was achieved largely via the use of networks structures—which have already been proven to be inherent in most interactive systems and environments. A methodology to abstract tasks, model problems and environments as network graphs, and measure their complexities, was proposed. High-level directions for comparing different agent types operating in vastly different environments were given along with concrete examples. Finally, it was shown how one can connect the studies of intelligence to other spheres, notably business decision-making and management. Hence, this chapter accomplishes my third and last goal (**G03**) outlined in the introduction of this thesis.

In my future work I will elaborate on the methodology described in this chapter and apply it in order to analyse and compare intelligence among groups implementing several interaction strategies, (business) organisational structures and management practices similar to the ones depicted in Figure 6.11. The endeavour is to understand where the examined group behaviours overlap, which characteristics they share, and what is peculiar to each sort of group behaviour or strategy examined.

Chapter 7

Conclusion and Future Work

A truly open mind means forcing our imaginations to conform to the evidence of reality, and not vice versa, whether or not we like the implications.

—Lawrence M. Krauss, *A Universe from Nothing* (2012)

By harnessing collective intelligence, more than 200,000 people from all around the world contributed to the classification of more than a hundred million galaxies. These classifications were in turn used to boost more than fifty research projects (Raddick et al., 2010). As of 2016, input received from the general public will decide what areas NASA’s Juno Mission will observe on Jupiter (Hansen et al., 2012). By querying people on the web, mobile applications such as VizWiz can answer users’ (visual and audio) questions in nearly real-time (Bigham et al., 2010). Gamers collaborating and competing in online multi-player games have been successfully recruited to create accurate protein structure models (Khatib et al., 2011). These types of collaborative solutions known as online citizen science projects have solved problems in a wide range of scientific domains as discussed in (Bonney et al., 2009).

These and many other examples of collective intelligence are all around us. Nature, evolution and animal group behaviours have inspired the design of numerous meta-heuristics that contributed to solving complex problems and enhanced the design of AI systems and robotics. Due to principles like the wisdom of the crowds (Surowiecki, 2005) and others, human crowds have become forecast tools by polling their collective opinion and harnessing their collective knowledge using web-based applications such as Amazon Mechanical Turk, Wikipedia and TripAdvisor. Advantages of collective intelligence have spread across many disciplines like economics, law, psychology, sociology and politics. Collective intelligence has clearly transformed our daily business operations and activities. Hence, the more we understand collective intelligence and its characteristics, the more we can exploit its potential and reap its benefits, regardless of the type of the system in which it occurs.

7.1 Overall Contribution

Overall, the outcomes from this thesis have several implications on real-world problems from the design of intelligent multiagent systems, to the understanding of their social behaviour, and the prediction of their capacity for intelligence.

The conclusions drawn from the different chapters serve as general guidelines that give insight into how to explore the potential of collectives across different cognitive systems and research disciplines. For instance, these conclusions can be used to reveal the different configurations that maximise agent group performance, and to identify scenarios resulting in interesting collective phenomena such as the wisdom of the crowd (Surowiecki, 2005).

The thesis also provides experimental and theoretical results (partly from intelligence tests) which can be implemented to enhance business decision-making and management strategies by exploring their underlying group dynamics and organisational structures. Such results are in turn useful for improving the accuracy of performing analytical reasoning for some estimate and prediction-like problems, as for example those required in the fields of deliberations analytics (e.g., aggregating the contributions of participants), sports business-intelligence (e.g., predicting team performance, supply chain optimisation), consumer models (factors influencing consumer behaviour), social (online) commerce, collaborative filtering and recommendation systems (used in search engines like Google, and entertainment providers such as Netflix), social media and interaction (e.g., user-generated content in online applications such as LinkedIn and Facebook), etc.

In addition to the above, this thesis provides some new understanding of (human and non-human animal, robotics and swarm-like) agent-based systems by simulating their behaviours and showing how different types of communication modes and other parameters, as well as their inter-dependencies, influence the performance of such systems over search (e.g., foraging, predator/prey scenarios) and learning (e.g., inductive inference and compression) problems.

The remarkable similarities highlighted—after showing the existence of principles like modularity, adaptability and communication—among the different occurrences of collective intelligence make its nature more intriguing. Until now, intelligence has merely been explained as what is measured by an intelligence test. Whereas, this thesis presents a step forward towards understanding and measuring (collective) intelligence across different research disciplines by identifying its main underlying building-blocks and measuring their (individual and simultaneous) impact at the collective level. As this thesis identifies a unified set of rules responsible for the manifestation of collective intelligence in various cognitive systems, it will likely allow us to understand how these systems relate to one another and to connect the missing links between the wide range of relevant research areas.

7.2 Brief Summary and Main Outcomes

The thesis began with a review of the latest research conducted on collective intelligence and its measurement. Numerous studies were collected from different areas apropos: swarm intelligence and animal group behaviour, ecology, evolution and natural sciences,

artificial intelligence and multiagent systems, behavioural economics, social sciences and psychology, crowd-sourcing and human computation, complex adaptive systems, and many others, all in connection with the notion of collective intelligence. Intersections between these areas were discussed as for example, the influence of biology and evolution on artificial intelligence and complex problem solving (Section 2.2.1), as well as the impact of human groups and (online) communities on economics and markets (Section 2.2.2). Moreover, some properties and characteristics that relate to the behaviour of human and non-human collectives were described. In addition to presenting examples of collective intelligence in humans, animals and machines, this thesis reviewed in Section 2.3 some techniques and methodologies that have been used in the literature to evaluate and measure (individual and collective) intelligence in these systems and understand their behaviour.

I summarise the outcomes arising from this thesis in the paragraphs below.

In Chapter 3 I give a new understanding about agent group performance and its underlying dynamics. This is achieved by showing how a range of **factors and properties** that are inherent to agent groups of various cognitive types quantitatively **shape** the agents' **performances** across environments and problems of well-defined algorithmic information-theoretic complexities. For instance, I quantitatively measure the impact of a list of factors such as *task algorithmic information-theoretic complexity*, *the interaction mode between agents*, *their organisational structure*, *their observation and communication abilities and their decision-making dynamics*, as well as others, on agent individual and group effectiveness. By altering test settings such as the evaluation time of the agents or the *uncertainty* of the assessment problem, I show that it is possible for one cooperative system to outperform another under some values of these factors, although it fails to do so under others. I present results from intelligence tests showing that the imitation of expert agents does not always guarantee the group's optimal performance over these tests. Moreover, I measure the effect of introducing more agents into the group and show that their performance is tightly controlled by the communication protocol used between the group members. Not all collectives outperform their equally sized group of isolated agents. This is due to important factors such as *environment uncertainty* and the *algorithmic information-theoretic complexity of the assessment task*, which can control the capacity for intelligence in these groups. An important property of groups is their organisational structure. This property is very relevant to the efficiency of large organisations. In Chapter 3, I quantitatively show how collective agent performances significantly vary from one group organisational/network structure to another. In addition to the above outcomes, intrinsic dependencies between the examined factors are identified and measured in Chapter 4. Agents with low observation or perception abilities can be compensated for, and significantly improved, by increasing their communication entropies, thus leading to smarter systems. Nevertheless, introducing more communication between agents does not monotonically improve performance as commonly presumed. Overall, Chapters 3 and 4 investigate important factors that have been overlooked by many of the existing works on collective intelligence, and consequently give a new insight into their quantitative influence on group intelligence.

In Chapter 5 I adopt a more formal approach to the analysis and prediction of intelligence. A **new information-theoretic mathematical predictive model is devised** to predict the accuracy of agents over tasks of well-defined and quantifiable complexities. The model combines notions from a measurement paradigm in psychometrics, called Item Response Theory (IRT), and algorithmic-information theoretic intelligence tests. One major contribution from the proposed model is that it makes it possible to avert the perpetual need to simulate agents over intelligence tests every time we need to predict their performance under a different problem configuration. One can further infer from the model a lower bound on agent accuracy with respect to task complexity and the breadth of its solution space. Subsequently, the relationship between agent selection cost, task difficulty and accuracy can be formalised as an optimisation problem. The model allows us to predict the accuracy of various kinds of agents of different abilities over distinct problem settings. Important relationships between the accuracy of an agent, the complexity of the assessment task and the depth of its solution space are discussed. It is also shown how to use the model to analytically reason about the accuracy of agent collectives—something which cannot be disclosed from standard intelligence tests. On one hand, some given scenarios illustrate that one group of agents is more accurate than individual agents or other groups. On the other hand, other scenarios and settings reveal that the opposite is sometimes true.

Chapter 6 gives insights on how to practically and **quantitatively compare intelligence between** non-uniform types of agents, operating in **vastly different environments**. The premise of devising universal intelligence tests as general models used for this purpose is precluded after I show that the majority of intelligence tests are highly adaptive to the evaluated entity type or the problem/environment type. Thus an alternative approach is sought. Can network science be one of the media towards achieving such a goal? A network structure is ubiquitous in every model used to represent complex systems, and it is inherent in the state space of every automaton or array of automata (Green, 1994, 2000; Green and Bransden, 2006; Green, 2011). Thus, the majority of natural and artificial environments from biology, sociology, economics and artificial intelligence do have some underlying network structure. After looking into the relationships between (natural) environments, meaningful assessment tasks, network science and complexity, Chapter 6 describes how abstracting the nature of the studied entities and their environments as networks can be used as a tool to cross-fertilise a wide array of fields—despite their disparity. It also solicits the idea that coping with complexity is a genuine goal born with most natural or man-made systems, and that studying complex systems (Newman, 2011) to understand and measure collective intelligence is a promising direction considering its ample fields of application. In the context of these discussions, Chapter 6 presents a new formalisation for abstracting assessment tasks and modelling environments using network graphs and shows how to measure their complexities. This is further used as part of a preliminary methodology for comparing intelligence between different agent types including those operating in vastly distinct environments. Limitations and difficulties associated

with the proposed methodology are also outlined. Finally, studies of intelligence are connected to other spheres, notably business decision-making and management. It is shown how the work presented in this thesis can be used to enhance businesses by applying intelligence studies to analyse (their collective performances resulting from) their different organisational structures and managerial strategies.

7.3 Limitations and Directions for Future Work

Evaluation environments and settings: The intelligence tests used in Chapter 3 for the purpose of this research assess a set of important skills and abilities relevant to intelligence. Nonetheless, these skills and abilities do not account for the complete range of multiagent problems. An interesting research direction to further explore here is the measurement of coordination from agent interaction—an important feature in multiagent systems that can have significant influence on their performance. Problems that particularly require coordination can be designed and presented (using knowledge representation techniques) to evaluate multiagent systems in such a way that positive payoff occurs only if two or more agents perform some set of coordinated (possibly hierarchical) actions simultaneously. Similar settings can further be used to investigate how coordination might occur via exclusive ad-hoc interactions between agents. Such coordination scenarios might occur among human, animal or artificial agents.

Furthermore, multiagent settings in which agents *share* and *delegate* rewards are also worth exploring as they are basically two fundamental models of how agents (human animal or artificial) are motivated to collaborate. More sophisticated evaluation settings are those in which the (attributes of the) testing environments are dynamically updated and altered in response to the agent’s actions. A reactive evaluation environment can be used to assess how quickly a collective of agents can adapt to underlying real-time changes as opposed to individual agents working in isolation and other collectives.

Competition (between agents) is also an important factor that can be analysed to understand the collective behaviour of agent groups. For instance, competition can lead to better performing groups by providing additional motivation or incentives to the agents. This can be evaluated in environments where different rewards (or reward weights) are allocated according to the order or precedence of actions taken, or even in a setting where agents *eat* the rewards provided by the environment.

Overall, there exist many properties that contribute to the intelligence of (groups of) agents shared among different cognitive systems. This thesis looks at few of those—particularly, the ability of learning, doing inductive-inference, compression and search. In the future, I will extend my work in order to assess agents interacting over a larger set of environments, such that I can evaluate other types of properties that have not been addressed within the scope of this thesis.

Accuracy prediction and modelling: With regard to accuracy prediction modelling, one of my aims is to extend the model described in Chapter 5 in the future to address important topics in this area regarding label noise, classifier dominances, “rough sets”

and approximations of decision classes. This can significantly enhance our understanding of questions relevant to task hardness, noise handling, outliers, meta-learning, etc. A more rigorous approximation of the true error will also be sought using a variant of *PAC Learning* theory (Valiant, 1984). Moreover, applications over complex and sophisticated tasks requiring a range of abilities (e.g., logics and spatial orientation) for success can be investigated. With respect to group accuracies, more sophisticated voting rules will be used in the future to analytically reason about team accuracy by analysing the outcomes from different sampling techniques over the agents' ranked votes. This is likely to have strong implications on a wide range of research disciplines (e.g., health, finance and forecasting services) where machine learning classification plays an important role, as well as other areas of AI such as game-playing and machine intelligence testing.

As mentioned in Chapter 5, the measurement of agent accuracies assumes that there is a selection of discrete *solutions* to choose from. While this reflects some important dimensions of intelligence linked to an agent's general ability (e.g., perform classification and inference tasks), most business solutions are in fact plans of action that list what needs to be done and how to divide up the work between team members. For example, such plans would take the form of constructed sets/sequences and appropriate task allocations to the group members that lead to specialisation. Therefore, the work in Chapter 5 could be extended to further construct accurate task allocation models for groups of heterogeneous agents, and link the results to realistic business operation scenarios.

Encoding environments and tasks: Chapter 6 proposes an approach for the quantitative comparison of performance between diverse kinds of agents operating in substantially different environments. This was made possible by providing common grounds for evaluation, which is in turn accomplished by presenting a methodology to model such environments as network graphs and measure their complexities prior to evaluation. However, many simplification and abstractions were made with respect to the encoding of the *states* of a problem/environment (e.g., labelling of the vertices in a network that abstract the problem/environment states). Hence, there is still much room for improvement in this regard. In addition, a promising future direction in this line of work could be to devise or find alternative, more elaborate and possibly composite measures of complexity that accurately quantify the difficulty of particular kinds of environment and their assessment tasks (especially tasks that are usually needed in businesses and whose complexity can be very difficult to measure). This can possibly still be achieved by modelling environments of interest as network graphs and consequently applying the new (information-theoretic) complexity measures to their corresponding graphs. Having elaborate general measures of (different types of problem/graph) complexity is likely to have strong implications on many research fields in which network graphs are meaningful especially those related to agent-based modelling.

Moreover, I intend to apply the methodology described in Chapter 6 in order to analyse and compare intelligence among groups implementing several interaction strategies, business organisational structures and management practices like the ones illustrated in

Figure 6.11. The endeavour is to understand where the examined group behaviours overlap, which characteristics they share, and what is peculiar to each sort of group behaviour or strategy examined.

Last Words

At the end this thesis, I hope that I have persuaded the readers of the importance of quantitatively investigating collective intelligence across the different cognitive systems human, animal and machine. There is plenty to be done in the future and there exist many clear avenues for further work and improvement in this area of research. We have seen that the wide range of research studies on collective intelligence really address somewhat similar questions, albeit their distinctiveness. Put in a nutshell, these questions look into the characteristics shaping the intelligence of agent-based systems and allowing them to cope with complex environments and solve difficult tasks. As a result, the process of understanding collective intelligence cannot be confined to a particular type of agents, nor to a single research discipline. It is rather the case that the wide range of studies exploring this interesting phenomenon complement and learn from one another. Moreover, this goal cannot be completely achieved without combining quantitative measurement, experimentation and (predictive) modelling, all of which play a central role in understanding the dynamics of interactive agent systems and predicting their collective performance and behaviour.

In recent years, the area of Artificial General Intelligence (AGI) has made important breakthroughs in designing, developing and testing (general-purpose) intelligent artificial agents. This will perhaps lead to scenarios where artificial agents cannot be distinguished from human agents in the context of intelligence testing in the future. For example, (Kirkpatrick et al., 2017) recently proposed a solution to one of the major problems in AI, namely the *catastrophic forgetting* in neural networks. Results from (Kirkpatrick et al., 2017) enable artificial models to learn—for the first time—multiple tasks sequentially, just like humans do. Moreover, formal measures of intelligence have been provided for artificial agents that are fully embedded within their environment (Orseau and Ring, 2012). This is interesting since agents are not simply studied in a setting where they interact with an external world, but rather their intelligence is “computed by, can be modified by, and is subject to the time and space constraints of the environment” with which they interacts (Orseau and Ring, 2012). These types of research achievements will likely lead to a major shift in how we understand and measure intelligence in artificial systems, and possibly human systems as well.

Advances in research on deep reinforcement-learning and deep neural networks run by the *Deepmind* company has pushed the boundaries of AI, especially with the first computer program to ever beat a professional player at the game of “Go” (Silver et al., 2016) and the better-than-human-level performing agents on classic Atari games (Schölkopf, 2015). As we get closer and closer to super-intelligent AI, many future risks can possibly arise from intelligent agents. These risks can potentially present dire challenges to humans as pointed out by (Bostrom, 2014). Examples include artificial agents misidentifying their objectives,

resistance to change goals, *instrumentality* (e.g., humans as resources and the paperclip AI problem) and more importantly *unpredictability* where intelligent AI agents cannot be certified for unseen situations. Solomonoff also discussed the dangers of a very intelligent machine (Solomonoff, 1967). AI risks are not confined to the future. Financial disasters have already occurred due to unpredictable artificial agent behaviours as for example “the 2010 Crash of 2:45” (Kirilenko et al., 2017) which was described as one of the most turbulent periods in the history of financial markets. Some undesirable consequences of endowing an intelligent agent with the ability to modify its own code have been considered in (Ring and Orseau, 2011). Moreover, prominent intellectuals such as Stephen Hawking and those affiliated with the likes of the Future of Humanity and the Future of Life institutes have repeatedly warned about concrete and prominent problems of AI such as machine learning-based systems controlling industrial processes, health-related systems, and other mission-critical technology (Amodei et al., 2016). It is still unclear whether many of these issues originate from too little or too much machine/artificial intelligence, or as a result of some emergent phenomena from the interaction of more than one agent.

Given the above considerations, novel (collective) intelligence testing frameworks might need to be devised in order to assess and predict AI risks and their impact. Moreover, the design of intelligence tests should become much more sophisticated in the endeavour to accommodate for the wide range of potential intelligence levels, including super machine-intelligence, although it is not clear at the time being whether or not it is feasible for humans to design such tests.

References

- Abbass, H. A. (2001). MBO: Marriage in honey bees optimization - A haplometrosis polygynous swarming approach, *Evolutionary Computation, 2001. Proceedings of the 2001 Congress on*, Vol. 1, IEEE, pp. 207–214.
- Adams, S., Arel, I., Bach, J., Coop, R., Furlan, R., Goertzel, B., Hall, J. S., Samsonovich, A., Scheutz, M., Schlesinger, M. et al. (2012). Mapping the landscape of human-level artificial general intelligence, *AI Magazine* **33**(1): 25–42, Association for the Advancement of Artificial Intelligence.
- Aggarwal, C. C. and Han, J. (2014). *Frequent pattern mining*, Springer.
- Aggarwal, I. and Woolley, A. W. (2013). Do you see what I see? The effect of members' cognitive styles on team processes and errors in task execution, *Organizational Behavior and Human Decision Processes* **122**(1): 92–99, Elsevier.
- Agusta, Y. and Dowe, D. L. (2002). MML clustering of continuous-valued data using Gaussian and t distributions, in B. McKay and J. Slaney (eds), *Proceedings of the 15th Australian Joint Conference on Artificial Intelligence*, Vol. 2557 of *Lecture Notes in Artificial Intelligence (LNAI)*, Springer-Verlag, Berlin, Germany, pp. 143–154.
- Agusta, Y. and Dowe, D. L. (2003a). Unsupervised learning of correlated multivariate Gaussian mixture models using MML, in T. D. Gedeon and L. C. Fung (eds), *Proceedings of the 16th Australasian Joint Conference on Artificial Intelligence*, Vol. 2903 of *Lecture Notes in Artificial Intelligence*, Springer-Verlag, Berlin, Germany, pp. 477–489.
- Agusta, Y. and Dowe, D. L. (2003b). Unsupervised learning of gamma mixture models using minimum message length, in M. H. Hamza (ed.), *Proceedings of the 3rd IASTED Conference on Artificial Intelligence and Applications*, ACTA Press, Benalmadena, Spain, pp. 457–462.
- Aliu, O. G., Imran, A., Imran, M. A. and Evans, B. (2013). A survey of self organisation in future cellular networks, *Communications Surveys & Tutorials* **15**(1): 336–361, IEEE.
- Aluru, K., Tellex, S., Oberlin, J. and MacGlashan, J. (2015). Minecraft as an experimental world for AI in robotics, *2015 AAAI Fall Symposium Series*.
- Amaneddine, N. and Chmait, N. (2009). Modeling real-estate search service using GIS technology, *17th International Conference on Geoinformatics*, IEEE, Fairfax, Virginia,

- USA, pp. 1–5.
URL: <http://ieeexplore.ieee.org/document/5293398/>
- Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J. and Mané, D. (2016). Concrete problems in AI safety, *arXiv preprint arXiv:1606.06565*.
URL: <https://arxiv.org/pdf/1606.06565.pdf>
- Angeline, P. J. (1995). Adaptive and self-adaptive evolutionary computations, *Computational intelligence: a dynamic systems perspective*, Citeseer.
- Anicich, E. M., Swaab, R. I. and Galinsky, A. D. (2015). Hierarchical cultural values predict success and mortality in high-stakes teams, *Proceedings of the National Academy of Sciences* **112**(5): 1338–1343, National Academy of Sciences.
- Ariely, D. (2010). *The upside of irrationality: The Unexpected Benefits of Defying Logic at Work and at Home*, Harper.
- Bachrach, Y., Graepel, T., Kasneci, G., Kosinski, M. and Van Gael, J. (2012). Crowd IQ: aggregating opinions to boost performance, *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems-Volume 1*, International Foundation for Autonomous Agents and Multiagent Systems, pp. 535–542.
- Banerjee, A. V. (1992). A simple model of herd behavior, *The Quarterly Journal of Economics* pp. 797–817, JSTOR.
- Barabási, A.-L. and Albert, R. (1999). Emergence of scaling in random networks, *Science* **286**(5439): 509–512, American Association for the Advancement of Science.
- Barillot, E., Calzone, L., Hupe, P., Vert, J.-P. and Zinovyev, A. (2012). *Computational systems biology of cancer*, CRC Press.
- Barmpalias, G. and Dowe, D. L. (2012). Universality probability of a prefix-free machine, *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences* **370**(1971): 3488–3511, The Royal Society. Theme Issue “The foundations of computation, physics and mentality: the Turing legacy” compiled and edited by Barry Cooper and Samson Abramsky.
- Barnette, D. W. (1983). *Map coloring, polyhedra, and the four-color problem*, number 8 in *The Dolciani Mathematical Expositions*, Mathematical Association of America.
- Beattie, C., Leibo, J. Z., Teplyashin, D., Ward, T., Wainwright, M., Küttler, H., Lefrancq, A., Green, S., Valdés, V., Sadik, A. et al. (2016). Deepmind lab, *arXiv preprint arXiv:1612.03801*.
URL: <https://arxiv.org/pdf/1612.03801.pdf>
- Bellmore, M. and Nemhauser, G. L. (1968). The traveling salesman problem: a survey, *Operations Research* **16**(3): 538–558, INFORMS.

- Beni, G. (2004). From swarm intelligence to swarm robotics, *Swarm robotics*, Springer, pp. 1–9.
- Benkler, Y. (2006). *The wealth of networks: How social production transforms markets and freedom*, Yale University Press.
- Bennett, G. K., Seashore, H. G. and Wesman, A. G. (1956). The differential aptitude tests: An overview, *The Personnel and Guidance Journal* **35**(2): 81–91, Wiley Online Library.
- Berinsky, A. J., Huber, G. A. and Lenz, G. S. (2012). Evaluating online labor markets for experimental research: Amazon.com’s Mechanical Turk, *Political Analysis* **20**(3): 351–368, SPM-PMSAPSA.
- Bernon, C., Chevrier, V., Hilaire, V. and Marrow, P. (2006). Applications of self-organising multi-agent systems: An initial framework for comparison, *Informatica (Slovenia)* **30**(1): 73–82.
- Bettencourt, L. M. A. (2009). The rules of information aggregation and emergence of collective intelligent behavior, *Topics in Cognitive Science* **1**(4): 598–620, Blackwell Publishing Ltd.
URL: <http://dx.doi.org/10.1111/j.1756-8765.2009.01047.x>
- Bien, Z., Bang, W.-C., Kim, D.-Y. and Han, J.-S. (2002). Machine Intelligence Quotient: its measurements and applications, *Fuzzy Sets and Systems* **127**(1): 3–16.
URL: <http://www.sciencedirect.com/science/article/pii/S016501140100149X>
- Bigham, J. P., Jayant, C., Ji, H., Little, G., Miller, A., Miller, R. C., Miller, R., Tatarowicz, A., White, B., White, S. et al. (2010). Vizwiz: nearly real-time answers to visual questions, *Proceedings of the 23rd annual ACM symposium on User interface software and technology*, ACM, pp. 333–342.
- Bikhchandani, S., Hirshleifer, D. and Welch, I. (1992). A theory of fads, fashion, custom, and cultural change as informational cascades, *Journal of Political Economy* **100**(5): 992–1026, JSTOR.
- Binet, A. and Simon, T. (1904). Méthodes nouvelles pour le diagnostic du niveau intellectuel des anormaux, *L’année Psychologique* **11**(1): 191–244, Persée-Portail des Revues Scientifiques en SHS.
- Bird, C. D. and Emery, N. J. (2009). Rooks use stones to raise the water level to reach a floating worm, *Current Biology* **19**(16): 1410–1414, Elsevier.
- Birnbaum, A. (1968). Some latent trait models and their use in inferring an examinee’s ability, *Statistical theories of mental test scores* pp. 395–479, Addison-Wesley.
- Blum, L. and Blum, M. (1975). Toward a mathematical theory of inductive inference, *Information and Control* **28**(2): 125 – 155.
URL: <http://www.sciencedirect.com/science/article/pii/S0019995875902612>

- Bonabeau, E. (2009). Decisions 2.0: The power of collective intelligence, *MIT Sloan Management Review* **50**(2): 45–52, Massachusetts Institute of Technology.
- Bonabeau, E., Dorigo, M. and Theraulaz, G. (1999). *Swarm Intelligence: From Natural to Artificial Systems*, Oxford University Press, Inc., New York, NY, USA.
- Bonabeau, E. and Meyer, C. (2001). Swarm Intelligence: A Whole New Way to Think About Business, *Harvard Business Review* **79**(5): 106–115.
URL: http://www.pina.ch/download/HarvardBR_SwarmIntelligence.pdf
- Bonchev, D. (1995). Kolmogorov’s information, Shannon’s entropy, and topological complexity of molecules, *Bulgarian Chemical Communications* **28**(3-4): 567–582, Bulgarian Academy of Sciences.
- Bonchev, D. (2004). Complexity analysis of yeast proteome network, *Chemistry & Biodiversity* **1**(2): 312–326, Wiley Online Library.
- Bonney, R., Cooper, C. B., Dickinson, J., Kelling, S., Phillips, T., Rosenberg, K. V. and Shirk, J. (2009). Citizen science: a developing tool for expanding science knowledge and scientific literacy, *BioScience* **59**(11): 977–984, BioOne.
- Bosse, T., Jonker, C. M., Schut, M. C. and Treur, J. (2006). Collective representational content for shared extended mind, *Cognitive Systems Research* **7**(2-3): 151–174, Elsevier Science Publishers.
URL: <http://dx.doi.org/10.1016/j.cogsys.2005.11.007>
- Bostrom, N. (2014). *Superintelligence: Paths, dangers, strategies*, Oxford University Press, Oxford, UK.
- Botoeva, E. (2012). Poppy flower TikZ script, TExample.net repository (Online). Last-accessed on February, 2017.
URL: <http://www.texample.net/tikz/examples/poppy/>
- Boulton, D. M. (1970). Numerical classification based on information measure, Master’s thesis, Basser Computing Department, University of Sydney, Sydney, Australia.
- Boulton, D. M. (1975). *The information measure criterion for intrinsic classification*, PhD thesis, Department of Computer Science, Monash University, Clayton, Australia.
- Boulton, D. M. and Wallace, C. S. (1969). The information content of a multistate distribution, *Journal of Theoretical Biology* **23**(2): 269–278, Elsevier.
- Boulton, D. M. and Wallace, C. S. (1970). A program for numerical classification, *The Computer Journal* **13**(1): 63–69, British Computer Society.
- Boulton, D. M. and Wallace, C. S. (1973a). A comparison between information measure classification, *Proceedings of the Australian & New Zealand Association for the Advancement of Science (ANZAAS) Congress (abstract)*, Perth, Australia.

- Boulton, D. M. and Wallace, C. S. (1973b). An information measure for hierarchic classification, *The Computer Journal* **16**(3): 254–261, British Computer Society.
- Boulton, D. M. and Wallace, C. S. (1973c). Occupancy of a rectangular array, *The Computer Journal* **16**(1): 57–63, British Computer Society.
- Boulton, D. M. and Wallace, C. S. (1975). An information measure for single link classification, *The Computer Journal* **18**(3): 236–238, British Computer Society.
- Bourjot, C., Chevrier, V. and Thomas, V. (2003). A new swarm mechanism based on social spiders colonies: from web weaving to region detection, *Web Intelligence and Agent Systems: An International Journal* **1**(1): 47–64, IOS Press.
- Brabham, D. C. (2010). Moving the crowd at Threadless: Motivations for participation in a crowdsourcing application, *Information, Communication & Society* **13**(8): 1122–1145, Taylor & Francis.
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J. and Zaremba, W. (2016). OpenAI gym, *arXiv preprint arXiv:1606.01540*.
URL: <https://arxiv.org/pdf/1606.01540.pdf>
- Brooks, R. A. (1991). Intelligence without reason, *Proceedings of the 1991 International Joint Conference on Artificial Intelligence*, Morgan Kaufmann, pp. 569–595.
URL: <http://www.ai.mit.edu/people/brooks/papers/AIM-1293.ps.Z>
- Brown, P. and Lauder, H. (2001). Human capital, social capital and collective intelligence, in S. Baron, J. Field and T. Schuller (eds), *Social Capital: Critical Perspectives*, Oxford University Press, Oxford, pp. 226–242.
URL: <http://opus.bath.ac.uk/17728/>
- Buhrmester, M., Kwang, T. and Gosling, S. D. (2011). Amazon’s Mechanical Turk a new source of inexpensive, yet high-quality, data?, *Perspectives on Psychological Science* **6**(1): 3–5, Sage Publications.
- Burbank, N., Dutta, D., Goel, A., Lee, D., Marschner, E. and Shivakumar, N. (2011). Widescope - a social platform for serious conversations on the web, *arXiv preprint arXiv:1111.1958*.
URL: <https://arxiv.org/pdf/1111.1958.pdf>
- Byrski, A. and Kisiel-Dorohinicki, M. (2003). Collective intelligence from a population of evolving neural networks, *Proceedings of the Intelligent Information Processing and Web Mining.*, Vol. 22 of *Advances in Soft Computing (AINSC)*, Springer, Zakopane, Poland, pp. 401–410.
- Cachia, R., Compañó, R. and Da Costa, O. (2007). Grasping the potential of online social networks for foresight, *Technological Forecasting and Social Change* **74**(8): 1179–1203, Elsevier.

- Camerer, C. F. and Fehr, E. (2006). When does “economic man” dominate social behavior?, *Science* **311**(5757): 47–52, American Association for the Advancement of Science.
- Camerer, C. F., Loewenstein, G. and Rabin, M. (2011). *Advances in behavioral economics*, Princeton University Press.
- Castelfranchi, C. (1998). Modelling social action for ai agents, *Artificial Intelligence* **103**(1-2): 157–182, Elsevier.
- Chaitin, G. J. (1966). On the length of programs for computing finite binary sequences, *Journal of the ACM (JACM)* **13**(4): 547–569, ACM.
- Chaitin, G. J. (1969). On the length of programs for computing finite binary sequences: statistical considerations, *Journal of the ACM (JACM)* **16**(1): 145–159, ACM.
- Chaitin, G. J. (1982). Godel’s theorem and information, *International Journal of Theoretical Physics* **21**(12): 941–954.
- Chaitin, G. J. (2002). On the intelligibility of the universe and the notions of simplicity, complexity and irreducibility, *arXiv preprint math/0210035*.
URL: <https://arxiv.org/pdf/math/0210035.pdf>
- Charness, G. and Sutter, M. (2012). Groups make better self-interested decisions, *The Journal of Economic Perspectives* **26**(3): 157–176, JSTOR.
- Cheeseman, P., Kanefsky, B. and Taylor, W. M. (1991). Where the really hard problems are, *Proceedings of the 12th International Joint Conference on Artificial Intelligence*, Vol. 91 of *IJCAI’91*, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, pp. 331–337.
URL: <http://dl.acm.org/citation.cfm?id=1631171.1631221>
- Child, J. (1972). Organizational structure, environment and performance: The role of strategic choice, *Sociology* **6**(1): 1–22, Sage Publications.
- Chmait, N. (2016). The Lambda Star intelligence test code-base, GitHub Repository.
URL: <https://github.com/nader-chmait/LambdaStar.git>
- Chmait, N. and Challita, K. (2013). Using simulated annealing and ant-colony optimization algorithms to solve the scheduling problem, *Computer Science and Information Technology* **1**(3): 208–224, Horizon Research Publishing, USA.
- Chmait, N., Dowe, D. L., Green, D. G. and Li, Y.-F. (2015). Observation, communication and intelligence in agent-based systems, in J. Bieger, B. Goertzel and A. Potapov (eds), *Proceedings of the 8th International Conference on Artificial General Intelligence*, Vol. 9205 of *Lecture Notes in Artificial Intelligence (LNAI)*, Springer, Berlin, Germany, pp. 50–59.
URL: http://dx.doi.org/10.1007/978-3-319-21365-1_6

- Chmait, N., Dowe, D. L., Li, Y.-F. and Green, D. G. (2017). An information-theoretic predictive model for the accuracy of AI agents adapted from psychometrics, *in* T. Everitt, B. Goertzel and A. Potapov (eds), *Proceedings of the 10th International Conference on Artificial General Intelligence (AGI 2017), August 15-18*, Vol. 10414 of *Lecture Notes in Artificial Intelligence (LNAI), Chap. 21*, Springer International Publishing, Melbourne, VIC, Australia, pp. 225–236. **Winner of the 2017 Kurzweil Best Paper Prize.**
URL: https://doi.org/10.1007/978-3-319-63703-7_21
- Chmait, N., Dowe, D. L., Li, Y.-F., Green, D. G. and Insa-Cabrera, J. (2015). Measuring universal intelligence in agent-based systems using the Anytime Intelligence Test, *Technical Report 2015/279*, Faculty of Information Technology (FIT), Clayton, Monash University, Australia.
URL: <http://www.csse.monash.edu.au/publications/2015/tr-2015-279-full.pdf>
- Chmait, N., Dowe, D. L., Li, Y.-F., Green, D. G. and Insa-Cabrera, J. (2016). Factors of collective intelligence: How smart are agent collectives?, *in* G. A. Kaminka, M. Fox, P. Bouquet, E. Hüllermeier, V. Dignum, F. Dignum and F. van Harmelen (eds), *Proceedings of 22nd European Conference on Artificial Intelligence ECAI*, Vol. 285 of *Frontiers in Artificial Intelligence and Applications*, IOS Press, The Hague, The Netherlands, pp. 542–550.
URL: <http://ebooks.iospress.nl/volumearticle/44798>
- Chmait, N., Hernández-Orallo, J., Martínez-Plumed, F., Strannegård, C. and Thórisson, K. R. (2017). Second workshop on Evaluating General-Purpose AI (EGPAI-17). In conjunction with the International Joint Conference on Artificial Intelligence IJCAI-17, Melbourne, Australia.
URL: <http://www.dsic.upv.es/~flip/EGPAI2017/>
- Chmait, N., Li, Y.-F., Dowe, D. L. and Green, D. G. (2016). A dynamic intelligence test framework for evaluating AI agents, *Proceedings of 1st International Workshop on Evaluating General-Purpose AI (EGPAI 2016), European Conference on Artificial Intelligence (ECAI 2016)*, The Hague, The Netherlands, pp. 1–8.
URL: <http://www.ecai2016.org/content/uploads/2016/08/W14-EGPAI-2016.pdf>
- Colman, A. M. (2013). *Game theory and its applications: In the social and biological sciences*, Psychology Press.
- Colman, A. M. (2014). *Game theory and experimental games: The study of strategic interaction*, Vol. 4 of *International Series in Experimental Social Psychology. Colman, Andrew M. (Eds)*, Pergamon.
- Comley, J. W. and Dowe, D. L. (2003). General Bayesian networks and asymmetric languages, *Proceedings of the 2nd Hawaii International Conference on Statistics and Related Fields*, pp. 1–18.

- Comley, J. W. and Dowe, D. L. (2005). Minimum Message Length and generalized Bayesian nets with asymmetric languages, *Advances in Minimum Description Length: Theory and Applications, Chapter 11*, In P. D. Grünwald and I. J. Myung and M. A., Pitt (Eds) pp. 265–294, MIT Press. Final camera ready copy was submitted in October 2003.
- Conradt, L. and Roper, T. J. (2005). Consensus decision making in animals, *Trends in Ecology & Evolution* **20**(8): 449–456, Elsevier.
- Console, L., Lombardi, I., Picardi, C. and Simeoni, R. (2011). Toward a social web of intelligent things, *AI Communications* **24**(3): 265–279, IOS Press.
- Conte, R. (2002). Agent-based modeling for understanding social intelligence, *Proceedings of the National Academy of Sciences* **99**(suppl 3): 7189–7190, National Academy of Sciences.
- Couzin, I. D. (2009). Collective cognition in animal groups, *Trends in Cognitive Sciences* **13**(1): 36–43.
- Cuevas, E., Cienfuegos, M., Zaldívar, D. and Pérez-Cisneros, M. (2013). A swarm optimization algorithm inspired in the behavior of the social-spider, *Expert Systems with Applications* **40**(16): 6374–6384, Elsevier.
- Daft, R. (2012). *Organization theory and design*, Nelson Education.
- Dalton, R. and Dalton, N. (2017). How to escape a maze according to maths, The Conversation, Australia. Published on January 26, 2017 10.57 pm AEDT.
URL: <https://theconversation.com/how-to-escape-a-maze-according-to-maths-71582>
- Damper, R. I. (2000). Editorial for the special issue on ‘emergent properties of complex systems’: Emergence and levels of abstraction, *International Journal of Systems Science* **31**(7): 811–818, Taylor & Francis.
URL: <http://dx.doi.org/10.1080/002077200406543>
- De Ayala, R. J. (2013). *The theory and practice of item response theory*, Guilford Publications.
- De Caritat Marquis De Condorcet, M. (1785). *Essai sur l’application de l’analyse à la probabilité des décisions rendues à la pluralité des voix*, L’imprimerie Royale.
URL: <https://books.google.co.id/books?id=MyIOAAAAQAAJ>
- De Castro, L. N. and Timmis, J. (2002). *Artificial immune systems: a new computational intelligence approach*, Springer Science & Business Media.
- Deaner, R. O., Isler, K., Burkart, J. and Van Schaik, C. (2007). Overall brain size, and not Encephalization Quotient, best predicts cognitive ability across non-human primates, *Brain, Behavior and Evolution* **70**(2): 115–124, Karger Publishers.

- Denoyer, L. and Gallinari, P. (2006). The Wikipedia XML corpus, *International Workshop of the Initiative for the Evaluation of XML Retrieval*, Springer, pp. 12–19.
- Dimitrakakis, C., Hernández-Orallo, J., Strannegård, C. and Thórisson, K. R. (2016). First workshop on Evaluating General-Purpose AI (EGPAI-16). Run in conjunction with the 22nd European Conference on Artificial Intelligence ECAI, The Hague, The Netherlands.
URL: <http://users.dsic.upv.es/~flip/EGPAI2016/>
- Dorigo, M., Birattari, M. and Stützle, T. (2006). Ant colony optimization, *Computational Intelligence Magazine, IEEE* **1**(4): 28–39, IEEE.
- Dorigo, M. and Gambardella, L. M. (1997). Ant colonies for the travelling salesman problem, *BioSystems* **43**(2): 73–81, Elsevier.
- Dorigo, M. and Stützle, T. (2009). Ant colony optimization: overview and recent advances, *Technical Report TR/IRIDIA/2009-013*, IRIDIA, Institut de Recherches Interdisciplinaires et de Développements en Intelligence Artificielle, Université Libre de Bruxelles, Belgium.
- D’Oswaldo, E. (2016). USA TikZ map script, GitHub Public Repository (Online). Last accessed on February, 2017.
URL: <https://gist.github.com/bordaigorl/fce575813ff943f47505>
- Dowe, D. L. (2008a). Foreword re C. S. Wallace, *The Computer Journal* **51**(5): 523–560. Christopher Stewart WALLACE (1933-2004) memorial special issue.
URL: <http://dx.doi.org/10.1093/comjnl/bxm117>
- Dowe, D. L. (2008b). Minimum Message Length and statistically consistent invariant (objective?) Bayesian probabilistic inference from (medical) “evidence”, *Social Epistemology* **22**(4): 433–460, Taylor & Francis.
- Dowe, D. L. (2011). MML, hybrid Bayesian network graphical models, statistical consistency, invariance and uniqueness, in P. S. Bandyopadhyay and M. R. Forster (ed.), *Handbook of the Philosophy of Science*, Vol. 7 of *Philosophy of Statistics*, Elsevier, pp. 901–982.
- Dowe, D. L. (2013). Introduction to Ray Solomonoff 85th memorial conference, in D. L. Dowe (ed.), *Algorithmic Probability and Friends. Bayesian Prediction and Artificial Intelligence*, Vol. 7070 of *Lecture Notes in Computer Science*, Springer, Berlin, Heidelberg, pp. 1–36.
URL: http://dx.doi.org/10.1007/978-3-642-44958-1_1
- Dowe, D. L. and Hajek, A. R. (1997a). A computational extension to the Turing Test, *Proceedings of the 4th Conference of the Australasian Cognitive Science Society, University of Newcastle, NSW, Australia*, Vol. 1, Citeseer.
URL: <http://users.monash.edu/~dld/Publications/1997/DoweHajek1997a.pdf>

- Dowe, D. L. and Hajek, A. R. (1997b). A computational extension to the Turing Test, *Technical Report #97/322*, Department of Computer Science, Monash University, Melbourne, Australia.
URL: <http://users.monash.edu/~dld/Publications/1997/DoweHajek1997b.pdf>
- Dowe, D. L. and Hajek, A. R. (1998). A non-behavioural, computational extension to the Turing Test, *International conference on computational intelligence & multimedia applications (ICCIMA'98), Gippsland, Australia*, pp. 101–106.
URL: <http://users.monash.edu/~dld/Publications/1998/DoweHajek1998.pdf>
- Dowe, D. L. and Hernández-Orallo, J. (2012). IQ tests are not for machines, yet, *Intelligence* **40**(2): 77–81, Elsevier.
URL: <http://users.monash.edu/~dld/Publications/2012/IQnotuniversal.pdf>
- Dowe, D. L. and Hernández-Orallo, J. (2014). How universal can an intelligence test be?, *Adaptive Behavior - Animals, Animats, Software Agents, Robots, Adaptive Systems* **22**(1): 51–69, Sage Publications, Inc.
URL: <http://dx.doi.org/10.1177/1059712313500502>
- Dowe, D. L., Hernández-Orallo, J. and Das, P. K. (2011). Compression and intelligence: Social environments and communication, *Proceedings of the 4th International Conference on Artificial General Intelligence*, Springer, Berlin, pp. 204–211.
URL: <http://dl.acm.org/citation.cfm?id=2032873.2032895>
- Dowe, D. L. and Korb, K. B. (1996). Conceptual difficulties with the efficient market hypothesis: Towards a naturalized economics, in D. L. D. et al. (ed.), *Proceedings of Information, Statistics and Induction in Science*, World Scientific, Melbourne, Australia, pp. 212–223.
- Dowe, D. L., Oliver, J. J. and Wallace, C. S. (1996). MML estimation of the parameters of the spherical Fisher distribution, in S. Arikawa and A. K. Sharma (eds), *Proceedings of the 7th International Workshop on Algorithmic Learning Theory (ALT'96)*, Vol. 1160 of *Lecture Notes in Artificial Intelligence (LNAI)*, Springer, pp. 213–227.
- Dowe, D. L. and Zaidi, N. A. (2010). Database normalization as a by-product of Minimum Message Length inference, *Proceedings of the 23rd Australasian Joint Conference on Artificial Intelligence*, Vol. 6464 of *Springer Lecture Notes in Artificial Intelligence (LNAI)*, Springer, Adelaide, Australia, pp. 82–91.
- Draves, S. (2008). Evolution and collective intelligence of the electric sheep, *The Art of Artificial Evolution*, Springer, pp. 63–78.
- Ducatelle, F., Di Caro, G. A. and Gambardella, L. M. (2010). Principles and applications of swarm intelligence for adaptive routing in telecommunications networks, *Swarm Intelligence* **4**(3): 173–198, Springer.
- Easley, D. and Kleinberg, J. (2010). *Networks, crowds, and markets: Reasoning about a highly connected world*, Cambridge University Press.

- Eberhart, R. C. and Shi, Y. (2001). Particle swarm optimization: developments, applications and resources, *Proceedings of the 2001 Congress on evolutionary computation*, Vol. 1, IEEE, pp. 81–86.
- Edgoose, T. and Allison, L. (1999). MML Markov classification of sequential data, *Statistics and Computing* **9**(4): 269–278, Springer.
- Edgoose, T., Allison, L. and Dowe, D. L. (1998). An MML classification of protein structure that knows about angles and sequence, *Pacific Symposium on Biocomputing*, Vol. 3, World Scientific Publishing, pp. 585–596.
- Edwards, R. T. and Dowe, D. L. (1998). Single factor analysis in MML mixture modelling, *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, Vol. 1394 of *Lecture Notes in Artificial Intelligence (LNAI)*, Springer, pp. 96–109.
- Eguíluz, V. M., Zimmermann, M. G., Cela-Conde, C. J. and San Miguel, M. (2005). Cooperation and the emergence of role differentiation in the dynamics of social networks, *American Journal of Sociology* **110**(4): 977–1008, JSTOR.
- Engel, D., Woolley, A. W., Aggarwal, I., Chabris, C. F., Takahashi, M., Nemoto, K., Kaiser, C., Kim, Y. J. and Malone, T. W. (2015). Collective intelligence in computer-mediated collaboration emerges in different contexts and cultures, *proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems CHI'15*, ACM, New York, NY, USA, pp. 3769–3778.
URL: <http://doi.acm.org/10.1145/2702123.2702259>
- Engelbrecht, A. P. (2006). *Fundamentals of computational swarm intelligence*, John Wiley & Sons.
- Eppstein, D. (1995). Subgraph isomorphism in planar graphs and related problems, *Proceedings of the Sixth Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA'95, Society for Industrial and Applied Mathematics, pp. 632–640.
URL: <http://dl.acm.org/citation.cfm?id=313651.313830>
- Evans, S., Hershey, J. and Saulnier, G. (2002). Kolmogorov complexity estimation and analysis, *6th World Conf. on Systemics, Cybernetics and Informatics, session on Complexity Theory and its applications to Systems, Networks and Information Assurance*. Technical Report 2002GRC177.
URL: <http://www.detectingdesign.com/PDF%20Files/Kolmogorov%20Complexity%202.pdf>
- Evans, T. G. (1964). A heuristic program to solve geometric-analogy problems, *Proceedings of the April 21-23, 1964, Spring Joint Computer Conference*, AFIPS '64 (Spring), ACM, pp. 327–338.
- Fallenstein, B. and Soares, N. (2014). Problems of self-reference in self-improving space-time embedded intelligence, in B. Goertzel, L. Orseau and J. Snider (eds), *Artificial*

- General Intelligence*, Vol. 8598 of *Lecture Notes in Computer Science*, Springer International Publishing, pp. 21–32.
URL: http://dx.doi.org/10.1007/978-3-319-09274-4_3
- Feldman, J. (2003). The Simplicity Principle in Human Concept Learning, *Current Directions in Psychological Science* **12**(6): 227–232.
URL: <http://dx.doi.org.monstera.cc.columbia.edu:2048/10.1046/j.0963-7214.2003.01267.x>
- Figueiredo, M. A. T. and Jain, A. K. (2002). Unsupervised learning of finite mixture models, *IEEE Transactions on pattern analysis and machine intelligence* **24**(3): 381–396, IEEE.
- Fiorito, G. and Scotto, P. (1992). Observational learning in octopus vulgaris, *Science* **256**(5056): 545–547, American Association for the Advancement of Science.
- Fitzgibbon, L. J., Dowe, D. L. and Vahid, F. (2004). Minimum message length autoregressive model order selection, in M. Palanaswami, C. C. Sekhar, G. K. Venayagamoorthy, S. Mohan and M. K. Ghantasala (eds), *Proceedings of International Conference on Intelligent Sensing and Information Processing*, IEEE, Chennai, India, pp. 439–444. Catalogue Number: 04EX783.
- Fong, T., Nourbakhsh, I. and Dautenhahn, K. (2003). A survey of socially interactive robots, *Robotics and Autonomous Systems* **42**(3): 143–166, Elsevier.
- Foster, K. R., Wenseleers, T. and Ratnieks, F. L. (2006). Kin selection is the key to altruism, *Trends in Ecology & Evolution* **21**(2): 57–60, Elsevier.
- Franklin, S. and Graesser, A. (1997). Is it an agent, or just a program?: A taxonomy for autonomous agents, in J. Mller, M. Wooldridge and N. Jennings (eds), *Intelligent Agents III Agent Theories, Architectures, and Languages*, Vol. 1193 of *Lecture Notes in Computer Science*, Springer Berlin Heidelberg, pp. 21–35.
URL: <http://dx.doi.org/10.1007/BFb0013570>
- Fulker, D. and Eysenck, M. (2012). *The Structure and Measurement of Intelligence*, Springer Berlin Heidelberg.
URL: <https://books.google.com.au/books?id=w0LVBQAAQBAJ>
- Galbraith, J. R. (2002). *Designing organizations. An executive guide to strategy, structure, and processes*, San Francisco: Jossey-Bass Publishers.
- Galef Jr, B. G. (1988). Imitation in animals: history, definition, and interpretation of data from the psychological laboratory, *Social Learning: Psychological and Biological Perspectives* **28**: 3–28, Lawrence Erlbaum Hillsdale.
- Galton, F. (1907). Vox populi (the wisdom of crowds), *Nature* **75**: 450–51.
- Gardner, R. A. and Gardner, B. T. (1969). Teaching sign language to a chimpanzee, *Science* **165**(3894): 664–672, Citeseer.

- Garnier, S., Gautrais, J. and Theraulaz, G. (2007). The biological principles of swarm intelligence, *Swarm Intelligence* **1**(1): 3–31, Springer.
- Gatti, D. D., Desiderio, S., Gaffeo, E., Cirillo, P. and Gallegati, M. (2011). *Macroeconomics from the Bottom-up*, Vol. 1, Springer Science & Business Media.
- Gatti, D. D., Di Guilmi, C., Gaffeo, E., Giulioni, G., Gallegati, M. and Palestrini, A. (2005). A new approach to business fluctuations: heterogeneous interacting agents, scaling laws and financial fragility, *Journal of Economic Behavior & Organization* **56**(4): 489–512, Elsevier.
- Gent, I. P. and Walsh, T. (1994). The SAT phase transition, *Proceedings of the 11th European Conference on Artificial Intelligence (ECAI), Amsterdam, the Netherlands*, Vol. 94, Pitman, pp. 105–109.
- Gibson, K. R., Rumbaugh, D. and Beran, M. (2001). Bigger is better: primate brain size in relationship to cognition, *Evolutionary Anatomy of the Primate Cerebral Cortex* pp. 79–97, Cambridge University Press.
- Gigerenzer, G. (2000). *Adaptive thinking: Rationality in the real world*, Oxford University Press, USA.
- Gigerenzer, G. and Goldstein, D. G. (1996). Reasoning the fast and frugal way: models of bounded rationality., *Psychological Review* **103**(4): 650, American Psychological Association.
- Gigerenzer, G. and Selten, R. (2002). *Bounded rationality: The adaptive toolbox*, MIT press.
- Gil, S. and Zanette, D. H. (2006). Coevolution of agents and networks: Opinion spreading and community disconnection, *Physics Letters A* **356**(2): 89–94, Elsevier.
- Gleizes, M.-P., Camps, V. and Glize, P. (1999). A theory of emergent computation based on cooperative self-organization for adaptive artificial systems, *Fourth European Congress of Systems Science*, pp. 20–24.
- Glenn, J. C., Gordon, T. J. and Florescu, E. (2014). *2013-14 State of the Future*, 1 edn, The Millennium Project.
- Goertzel, B. and Bugaj, S. V. (2009). AGI preschool: a framework for evaluating early-stage human-like AGIs, *Proceedings of the 2nd international conference on Artificial General Intelligence (AGI-09)*, Advances in Intelligent Systems Research, Atlantis Press, Arlington, Virginia, pp. 31–36.
- Goertzel, B., Iklé, M., Goertzel, I. F. and Heljakka, A. (2008). *Probabilistic logic networks: A comprehensive framework for uncertain inference*, Springer Science & Business Media.
- Gold, E. (1967). Language Identification in the Limit, *Information and Control* **10**(5): 447–474, Elsevier.

- Goldberg, D., Nichols, D., Oki, B. M. and Terry, D. (1992). Using collaborative filtering to weave an information tapestry, *Communications of the ACM* **35**(12): 61–70, ACM.
- Goldberg, L. R., Johnson, J. A., Eber, H. W., Hogan, R., Ashton, M. C., Cloninger, C. R. and Gough, H. G. (2006). The international personality item pool and the future of public-domain personality measures, *Journal of Research in Personality* **40**(1): 84–96, Elsevier.
- Goleman, D. (2007). *Social intelligence*, Random House.
- Gottfredson, L. S. (1997). Mainstream science on intelligence: An editorial with 52 signatories, history, and bibliography, *Intelligence* **24**(1): 13–23. First published in the Wall Street Journal, December 13, 1994.
- Grassé, P.-P. (1959). La reconstruction du nid et les coordinations interindividuelles chez *Bellicositermes natalensis* et *Cubitermes* sp. la théorie de la stigmergie: Essai d'interprétation du comportement des termites constructeurs, *Insectes Sociaux* **6**(1): 41–80, Springer.
- Gray, L. (2003). A mathematician looks at Wolfram's new kind of science, *Notices of the American Mathematical Society* **50**(2): 200–211.
URL: <http://www.ams.org/notices/200302/fea-gray.pdf>
- Green, D. G. (1994). Emergent behavior in biological systems, *Complexity International (electronic journal of complex systems research)* **1**: 1–12, Charles Stuart University.
URL: <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.111.1649&rep=rep1&type=pdf>
- Green, D. G. (2000). Self-organisation in complex systems, *Complex Systems* pp. 11–50, Cambridge University Press, Cambridge.
- Green, D. G. (2011). Elements of a network theory of complex adaptive systems, *International Journal of Bio-Inspired Computation* **13** **3**(3): 159–167, Inderscience Publishers Ltd.
- Green, D. G. (2014). *Of ants and men: The unexpected side effects of complexity in society*, Springer.
- Green, D. G. and Bransden, T. G. (2006). Complexity theory, *McGraw-Hill Encyclopedia of Science and Technology* pp. 507–511, McGraw-Hill, New York.
- Gregory, R. L. and Zangwill, O. L. (1987). *The Oxford companion to the mind*, Oxford University Press.
- Gruber, T. (2008). Collective knowledge systems: Where the social web meets the semantic web, *Web Semantics: science, services and agents on the World Wide Web* **6**(1): 4–13, Elsevier.

- Gubbi, J., Buyya, R., Marusic, S. and Palaniswami, M. (2013). Internet of things (IoT): A vision, architectural elements, and future directions, *Future Generation Computer Systems* **29**(7): 1645–1660, Elsevier.
- Guilford, J. P. (1967). *The nature of human intelligence*, McGraw-Hill.
- Guo, B., Zhang, D. and Wang, Z. (2011). Living with Internet of Things: The emergence of embedded intelligence, *Internet of Things (iThings/CPSCoM), 2011 International Conference on and 4th International Conference on Cyber, Physical and Social Computing*, IEEE, pp. 297–304.
- Gusfield, D. and Irving, R. W. (1989). *The stable marriage problem: structure and algorithms*, MIT press, Cambridge, MA, USA.
- Halmes, M. (2013). Measurements of collective machine intelligence, *arXiv preprint arXiv:1306.6649*. Master Thesis.
- Hamilton, W. D. (1964a). The genetical evolution of social behaviour. I, *Journal of Theoretical Biology* **7**(1): 1–16, Elsevier.
- Hamilton, W. D. (1964b). The genetical evolution of social behaviour. II, *Journal of Theoretical Biology* **7**(1): 17–52, Elsevier.
- Hansen, C., Bolton, S., Caplinger, M., Dyches, P., Jensen, E., Levin, S. and Ravine, M. (2012). JunoCam: A Public Endeavor, *AAS/Division for Planetary Sciences Meeting Abstracts*, Vol. 44, American Astronomical Society, DPS meeting #44, id.515.05.
- Hanus, D., Mendes, N., Tennie, C. and Call, J. (2011). Comparing the performances of apes (gorilla gorilla, pan troglodytes, pongo pygmaeus) and human children (homo sapiens) in the floating peanut task, *PloS One* **6**(6): e19555, Public Library of Science.
- Hardin, G. (1968). The tragedy of the commons, *Science* **162**(3859): pp. 1243–1248, American Association for the Advancement of Science.
URL: <http://www.jstor.org/stable/1724745>
- Hart, P. E., Nilsson, N. J. and Raphael, B. (1968). A formal basis for the heuristic determination of minimum cost paths, *IEEE Transactions on Systems Science and Cybernetics* **4**(2): 100–107, IEEE.
- Hashemi, V. and Endriss, U. (2014). Measuring diversity of preferences in a group, *21st European Conference on Artificial Intelligence ECAI 2014, Prague, Czech Republic - Including Prestigious Applications of Intelligent Systems (PAIS 2014)*, IOS Press, pp. 423–428.
- Havens, T. C., Spain, C. J., Salmon, N. G. and Keller, J. M. (2008). Roach infestation optimization, *Swarm Intelligence Symposium, 2008. SIS 2008. IEEE*, IEEE, pp. 1–7.
- Hearn, J. (2006). Applications of Kolmogorov complexity to graphs, Harvey Mudd College California, United States, Department of Mathematics, senior thesis.

- URL: <https://www.math.hmc.edu/seniorthesis/archives/2006/jhearn/jhearn-2006-thesis.pdf>
- Hernández-Orallo, J. (2000). Beyond the Turing Test, *Journal of Logic, Language and Information* **9**(4): 447–466, Springer.
- Hernández-Orallo, J. (2010). A (hopefully) unbiased universal environment class for measuring intelligence of biological and artificial systems, *Proceedings of 3rd Conference on Artificial General Intelligence*, AGI'10, Atlantis Press, pp. 182–183.
URL: <http://users.dsic.upv.es/proy/anynt/unbiased.pdf>
- Hernández-Orallo, J. (2015). On environment difficulty and discriminating power, *Autonomous Agents and Multi-Agent Systems* **29**(3): 402–454, Springer.
- Hernández-Orallo, J. (2017). *The measure of all minds: evaluating natural and artificial intelligence*, Cambridge University Press.
- Hernández-Orallo, J., Baroni, M., Bieger, J., Chmait, N., Dowe, D. L., Hofmann, K., Martínez-Plumed, F., Strannegård, C. and Thórisson, K. R. (2017). A New AI Evaluation Cosmos: Ready to Play the Game?, *AI Magazine, Association for the Advancement of Artificial Intelligence* **38**(3): 66–69.
URL: <https://www.aaai.org/ojs/index.php/aimagazine/article/view/2748>
- Hernández-Orallo, J. and Dowe, D. L. (2010). Measuring universal intelligence: Towards an anytime intelligence test, *Artificial Intelligence* **174**(18): 1508–1539, Elsevier Science Publishers Ltd.
URL: <http://dx.doi.org/10.1016/j.artint.2010.09.006>
- Hernández-Orallo, J. and Dowe, D. L. (2013). On potential cognitive abilities in the machine kingdom, *Minds and Machines* **23**(2): 179–210, Springer.
- Hernández-Orallo, J., Insa-Cabrera, J., Dowe, D. L. and Hibbard, B. (2012). Turing machines and recursive Turing tests, *AISB/IACAP 2012 Symposium “Revisiting Turing and his Test”*, Citeseer, pp. 28–33.
- Hernández-Orallo, J., Martínez-Plumed, F., Schmid, U., Siebers, M. and Dowe, D. L. (2016). Computer models solving intelligence test problems: Progress and implications, *Artificial Intelligence* **230**: 74–107, Elsevier.
URL: <http://www.sciencedirect.com/science/article/pii/S0004370215001538>
- Hernández-Orallo, J. and Minaya-Collado, N. (1998). A formal definition of intelligence based on an intensional variant of Kolmogorov complexity, *Proceedings of the Int. Symposium of EIS*, ICSC Press, pp. 146–163.
- Hillier, B. (2007). *Space is the machine: a configurational theory of architecture*, Space Syntax.
- Hingston, P. (2010). A new design for a Turing Test for bots, *2010 IEEE Symposium on Computational Intelligence and Games (CIG)*, IEEE, pp. 345–350.

- Hoffmann, M. and Pfeifer, R. (2012). The implications of embodiment for behavior and cognition: animal and robotic case studies, *arXiv preprint arXiv:1202.0440*.
URL: <https://arxiv.org/pdf/1202.0440>
- Hong, L. and Page, S. E. (2004). Groups of diverse problem solvers can outperform groups of high-ability problem solvers, *PNAS* **101**(46): 16385–16389.
URL: <http://www.pnas.org/content/101/46/16385.abstract>
- Howe, J. (2008). *Crowdsourcing: How the power of the crowd is driving the future of business*, Random House.
- Hutter, M. (2004). *Universal Artificial Intelligence: Sequential Decisions based on Algorithmic Probability*, EATCS, Springer, Berlin.
URL: <http://www.hutter1.net/ai/uaibook.htm>
- Insa-Cabrera, J., Benacloch-Ayuso, J.-L. and Hernández-Orallo, J. (2012). On measuring social intelligence: Experiments on competition and cooperation, in J. Bach, B. Goertzel and M. Iklé (eds), *Proceedings 5th International Conference on Artificial General Intelligence (AGI)*, Vol. 7716 of *Lecture Notes in Computer Science (LNCS)*, Springer Berlin Heidelberg, pp. 126–135.
URL: http://dx.doi.org/10.1007/978-3-642-35506-6_14
- Insa-Cabrera, J., Dowe, D. L., España-Cubillo, S., Hernández-Lloreda, M. V. and Hernández-Orallo, J. (2011). Comparing humans and AI agents, *Artificial General Intelligence (AGI)*, Vol. 6830 of *Lecture Notes in Computer Science (LNCS)*, Springer, pp. 122–132.
URL: <http://dblp.uni-trier.de/db/conf/agi/agi2011.html#Insa-CabreraDEHH11>
- Insa-Cabrera, J., Dowe, D. L. and Hernández-Orallo, J. (2011). Evaluating a reinforcement learning algorithm with a general intelligence test, *Conference of the Spanish Association for Artificial Intelligence*, Springer, pp. 1–11.
- Insa-Cabrera, J. and Hernández-Orallo, J. (2013). Interaction settings for measuring (social) intelligence in multi-agent systems, *ReteCog II Workshop: Interaction*, Zaragoza, Spain.
- Insa-Cabrera, J. and Hernández-Orallo, J. (2015). Instrumental properties of social testbeds, in J. Bieger, B. Goertzel and A. Potapov (eds), *Proceedings of the Eighth Conference on Artificial General Intelligence*, Vol. 9205 of *Lecture Notes in Artificial Intelligence (LNAI)*, Springer, pp. 101–110.
- Insa-Cabrera, J., Hernández-Orallo, J., Dowe, D. L., España, S. and Hernández-Lloreda, M. V. (2012). The ANYNT project intelligence test Lambda one, *AISB/IACAP 2012 Symposium “Revisiting Turing and his Test”*, pp. 20–27.
URL: <http://users.dsic.upv.es/~flip/papers/AISB-AICAP2012a.pdf>

- Ishii, H., Ogura, M., Kurisu, S., Komura, A., Takanishi, A., Iida, N. and Kimura, H. (2006). Experimental study on task teaching to real rats through interaction with a robotic rat, *Proceedings of the 9th International Conference on From Animals to Animats: Simulation of Adaptive Behavior*, Vol. 9 of *SAB'06*, Springer, pp. 643–654.
- Jackson, M. O. and Wolinsky, A. (2003). A strategic model of social and economic networks, *Networks and Groups*, Springer, pp. 23–49.
- Jerison, H. (2012). *Evolution of the Brain and Intelligence*, Elsevier.
- Jerison, H. J. (1975). *Evolution of the brain and intelligence*, Academic Press.
- Jiang, A., Marcolino, L. S., Procaccia, A. D., Sandholm, T., Shah, N. and Tambe, M. (2014). Diverse randomized agents vote to win, *Proceedings of the 27th International Conference on Neural Information Processing Systems*, NIPS'14, MIT Press, pp. 2573–2581.
- John, O. P., Naumann, L. P. and Soto, C. J. (2008). Paradigm shift to the integrative big five trait taxonomy, *Handbook of Personality: Theory and Research* **3**: 114–158, Guilford Press New York, NY.
- Johnson, M., Hofmann, K., Hutton, T. and Bignell, D. (2016). The Malmo platform for artificial intelligence experimentation, *International Joint Conference on Artificial Intelligence (IJCAI)*, p. 4246.
- Kao, A. B., Miller, N., Torney, C., Hartnett, A. and Couzin, I. D. (2014). Collective learning and optimal consensus decisions in social animal groups, *PLoS Computational Biology* **10**(8): e1003762, Public Library of Science.
- Karaboga, D. (2005). An idea based on honey bee swarm for numerical optimization, *Technical Report TR06*, Erciyes university, Engineering faculty, Computer Engineering Department, Kayseri, Türkiye.
URL: http://mf.erciyes.edu.tr/abc/pub/tr06_2005.pdf
- Karuna, H., Valckenaers, P., Saint-Germain, B., Verstraete, P., Zamfirescu, C. B. and Van Brussel, H. (2004). Emergent forecasting using a stigmergy approach in manufacturing coordination and control, *Engineering Self-Organising Systems*, Springer, pp. 210–226.
- Kasarapu, P. and Allison, L. (2015). Minimum message length estimation of mixtures of multivariate Gaussian and von Mises-Fisher distributions, *Machine Learning* **100**(2): 333–378, Springer.
- Kaspar, F. and Schuster, H. G. (1987). Easily calculable measure for the complexity of spatiotemporal patterns, *Physical Review A* **36**(2): 842–848, American Physical Society.
URL: <http://link.aps.org/abstract/PRA/v36/p842>
- Kasparov, G. K. and King, D. (2000). *Kasparov Against the World: The Story of the Greatest Online Challenge*, KasparovChess Online, Incorporated.

- Kauffman, S. A. (1993). *The origins of order: Self organization and selection in evolution*, Oxford University Press, USA.
- Kempka, M., Wydmuch, M., Runc, G., Toczek, J. and Jaśkowski, W. (2016). Vizdoom: A doom-based AI research platform for visual reinforcement learning, *arXiv preprint arXiv:1605.02097*.
URL: <https://arxiv.org/pdf/1605.02097.pdf>
- Kennedy, J. (2011). Particle swarm optimization, *Encyclopedia of machine learning*, Springer, pp. 760–766.
- Kennedy, J., Kennedy, J. F., Eberhart, R. C. and Shi, Y. (2001). *Swarm intelligence*, Morgan Kaufmann.
- Khatib, F., Cooper, S., Tyka, M. D., Xu, K., Makedon, I., Popović, Z. and Baker, D. (2011). Algorithm discovery by protein folding game players, *Proceedings of the National Academy of Sciences* **108**(47): 18949–18953, National Academy of Sciences.
- Kirilenko, A. A., Kyle, A. S., Samadi, M. and Tuzun, T. (2017). The flash crash: High frequency trading in an electronic market, *Journal of Finance* (Forthcoming). Available at SSRN: <https://ssrn.com/abstract=1686004>.
- Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., Hassabis, D., Clopath, C., Kumaran, D. and Hadsell, R. (2017). Overcoming catastrophic forgetting in neural networks, *Proceedings of the National Academy of Sciences*, published ahead of print March 14, 2017.
URL: <http://www.pnas.org/content/early/2017/03/13/1611835114.abstract>
- Kitano, H., Asada, M., Kuniyoshi, Y., Noda, I., Osawa, E. and Matsubara, H. (1997). Robocup: A challenge problem for AI, *AI Magazine* **18**(1): 73.
- Klein, M. and Garcia, A. C. B. (2015). Crowd computing: From human computation to collective intelligence (a tutorial), *Proceedings of 24th International Joint Conference on Artificial Intelligence IJCAI'2015*, Buenos Aires, Argentina.
URL: <http://ijcai-15.org/downloads/tutorials/T12-CollectiveIntelligence.pdf>
- Kolmogorov, A. N. (1965). Three approaches to the quantitative definition of information, *Problems of Information Transmission* **1**(1): 1–7.
- Koriat, A. (2012). When are two heads better than one and why?, *Science* **336**(6079): 360–362, American Association for the Advancement of Science.
- Kosinski, M., Bachrach, Y., Kasneci, G., Van-Gael, J. and Graepel, T. (2012). Crowd IQ: Measuring the intelligence of crowdsourcing platforms, *Proceedings of the 4th Annual ACM Web Science Conference*, ACM, pp. 151–160.

- Krause, S., James, R., Faria, J. J., Ruxton, G. D. and Krause, J. (2011). Swarm intelligence in humans: diversity can trump ability, *Animal Behaviour* **81**(5): 941–948, Elsevier.
- Kuncheva, L. I. (2004). *Combining Pattern Classifiers: Methods and Algorithms*, Wiley-Interscience.
- Laguna, G., Murrieta-Cid, R., Becerra, H. M., Lopez-Padilla, R. and LaValle, S. M. (2014). Exploration of an unknown environment with a differential drive disc robot, *2014 IEEE International Conference on Robotics and Automation (ICRA)*, IEEE, pp. 2527–2533.
- Landemore, H. (2013). *Democratic reason: Politics, collective intelligence, and the rule of the many*, Princeton University Press.
- Laughlin, P. R. (2011). *Group problem solving*, Princeton University Press.
- Legg, S. and Hutter, M. (2007). Universal intelligence: A definition of machine intelligence, *Minds and Machines* **17**(4): 391–444.
- Legg, S. and Veness, J. (2013). An approximation of the universal intelligence measure, *Algorithmic Probability and Friends. Bayesian Prediction and Artificial Intelligence*, Vol. 7070 of *Lecture Notes in Computer Science*, Springer, pp. 236–249.
- Leimeister, J. M. (2010). Collective intelligence, *Business & Information Systems Engineering* **2**(4): 245–248, Springer.
- Lempel, A. and Ziv, J. (1976). On the Complexity of Finite Sequences, *Information Theory, IEEE Transactions on* **22**(1): 75–81, IEEE.
URL: <http://dx.doi.org/10.1109/tit.1976.1055501>
- Levin, L. A. (1973). Universal sequential search problems, *Problems of Information Transmission* **9**(3): 265–266.
- Lévy, P. (1997). *Collective Intelligence: Mankind's Emerging World in Cyberspace*, Perseus Books, Cambridge, MA, USA.
- Li, M. and Vitányi, P. (2008). *An introduction to Kolmogorov complexity and its applications (3rd ed.)*, Springer-Verlag New York, Inc.
- Lihoreau, M., Raine, N. E., Reynolds, A. M., Stelzer, R. J., Lim, K. S., Smith, A. D., Osborne, J. L. and Chittka, L. (2012). Radar tracking and motion-sensitive cameras on flowers reveal the development of pollinator multi-destination routes over large spatial scales, *PLoS Biology* **10**(9): e1001392, Public Library of Science.
- Llopis, N. and Nicholson, C. (last accessed, April 2016). UnitTest++ testing framework, <https://github.com/unittest-cpp/unittest-cpp>.
- Lord, F. M. and Novick, M. R. (1968). *Statistical theories of mental test scores.*, Addison-Wesley.
URL: <http://ezproxy.library.yorku.ca/login?url=http://search.proquest.com/docview/615567141?accountid=15182>

- Lucas, S. M. (2007). Ms pac-man competition, *ACM SIGEVOlution* **2**(4): 37–38, ACM.
- Mačiulienė, M. (2014). Power through things: Following traces of collective intelligence in Internet of Things, *Socialnės Technologijos* **4**(1): 168–178, Mykolas Romeris University.
- Maes, P. (1993). Modeling adaptive autonomous agents, *Artificial Life* **1**(1_2): 135–162, MIT Press.
- Makalic, E. and Schmidt, D. F. (2011). MDL multiple hypothesis testing, *Proceedings of the 4th Workshop on Information Theoretic Methods in Science and Engineering (WITMSE-11)*, Citeseer, Helsinki, Finland, pp. 45–48. Invited paper.
- Malkiel, B. G. (1989). Efficient market hypothesis, *The New Palgrave: Finance*. Norton, New York pp. 127–134.
- Malone, T. W. (2008). What is collective intelligence and what will we do about it?, *Collective Intelligence: Creating a Prosperous World at Peace*, Earth Intelligence Network, M. Tovey (Ed.) pp. 1–4. Official launch of the MIT Center for Collective Intelligence (Academic Preface), Oakton, Virginia, Earth Intelligence Network.
- Malone, T. W. and Bernstein, M. S. (eds) (2015). *Handbook of Collective Intelligence*, The MIT Press.
URL: <https://mitpress.mit.edu/hci>
- Malone, T. W., Laubacher, R. and Dellarocas, C. (2010). The collective intelligence genome, *MIT Sloan Management Review* **51**(3): 21, Massachusetts Institute of Technology.
- Mamei, M., Vasirani, M. and Zambonelli, F. (2004). Self-organizing spatial shapes in mobile particles: The tota approach, *Engineering Self-Organising Systems*, Springer, pp. 138–153.
- March, J. G. (1991). Exploration and exploitation in organizational learning, *Organization Science* **2**(1): 71–87, INFORMS.
- Martinez-Plumed, F., Prudêncio, R. B., Martinez-Usó, A. and Hernández-Orallo, J. (2016). Making sense of item response theory in machine learning, *Proc. of 22nd Europ. Conf. on Artificial Int.*, Vol. 285 of *Frontiers in A. I. and Applications*, pp. 1140–1148.
- Mason, W. and Watts, D. J. (2012). Collaborative learning in networks, *PNAS* **109**(3): 764–769.
URL: <http://www.pnas.org/content/109/3/764.abstract>
- Matan, O. (1996). On voting ensembles of classifiers (extended abstract), *Proceedings of AAAI-96 workshop on integrating multiple learned models*, Citeseer, pp. 84–88.
- Maturana, F. P. and Norrie, D. H. (1996). Multi-agent mediator architecture for distributed manufacturing, *Journal of Intelligent Manufacturing* **7**(4): 257–270, Springer.

- May, K. O. (1952). A set of independent necessary and sufficient conditions for simple majority decision, *Econometrica: Journal of the Econometric Society* **20**(4): 680–684, JSTOR.
- Mengistu, H., Huizinga, J., Mouret, J.-B. and Clune, J. (2016). The evolutionary origins of hierarchy, *PLoS Computational Biology* **12**(6): e1004829, Public Library of Science.
- Michel, J.-B., Shen, Y. K., Aiden, A. P., Veres, A., Gray, M. K., Pickett, J. P., Hoiberg, D., Clancy, D., Norvig, P., Orwant, J. et al. (2011). Quantitative analysis of culture using millions of digitized books, *Science* **331**(6014): 176–182, American Association for the Advancement of Science.
- Mikolov, T., Joulin, A. and Baroni, M. (2015). A roadmap towards machine intelligence, *arXiv preprint arXiv:1511.08130*.
URL: <https://arxiv.org/pdf/1511.08130.pdf>
- Miller, N., Garnier, S., Hartnett, A. T. and Couzin, I. D. (2013). Both information and social cohesion determine collective decisions in animal groups, *Proceedings of the National Academy of Sciences* **110**(13): 5263–5268, National Academy of Sciences.
- Millonas, M. M. (1994). Swarms, phase transitions, and collective intelligence, *Proceedings of Santa Fe institute studies in the Sciences of complexity*, Vol. 17, Addison-Wesley publishing co, pp. 417–417.
- Mintzberg, H. (1979). *The structuring of organizations: A synthesis of the research*, Prentice-Hall.
- Mitchell, D., Selman, B. and Levesque, H. (1992). Hard and easy distributions of SAT problems, *Proceedings of the 10th National Conference on Artificial Intelligence, San Jose, California, Association for the Advancement of Artificial Intelligence (AAAI)*, Vol. 92, pp. 459–465.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G. et al. (2015). Human-level control through deep reinforcement learning, *Nature* **518**(7540): 529–533, Nature Research.
- Monismith, D. R. and Mayfield, B. E. (2008). Slime mold as a model for numerical optimization, *2008 IEEE Swarm Intelligence Symposium*, IEEE, pp. 1–8.
- Mowshowitz, A. and Dehmer, M. (2012). Entropy and the complexity of graphs revisited, *Entropy* **14**(3): 559–570, Molecular Diversity Preservation International.
- Musgrave, S. and Dowe, D. L. (2010). Kinship, optimality, and typology, *Behavioral and Brain Sciences (BBS)* **33**(5): 397–398, Cambridge University Press.
- Nagalakshmi, A. and Joglekar, S. (2011). Collective intelligence applications—algorithms and visualization, *Emerging Applications of Information Technology (EAIT), 2011 Second International Conference on*, IEEE, pp. 25–30.

- Namatame, A. and Chen, S.-H. (2016). *Agent Based Modelling and Network Dynamics*, Oxford University Press.
- Nehaniv, C. L. and Dautenhahn, K. (2007). *Imitation and social learning in robots, humans and animals: behavioural, social and communicative dimensions*, Cambridge University Press.
- Newman, M. E. (2006). Modularity and community structure in networks, *Proceedings of the National Academy of Sciences* **103**(23): 8577–8582, National Academy of Sciences.
- Newman, M. E. (2011). Complex systems: A survey, *American Journal of Physics* **79**(8): 800–810.
URL: <https://arxiv.org/pdf/1112.1440.pdf>
- Nolfi, S. and Floreano, D. (1998). Coevolving predator and prey robots: do arms races arise in artificial evolution?, *Artificial Life* **4**(4): 311–335, MIT Press.
- Noubel, J.-F. (2004). Collective intelligence, the invisible revolution. (the transitioner). Retrieved on 06 Jan 2017.
URL: http://www.oss.net/dynamaster/file_archive/070118/14da9d70ab635fb6f161a44fbf08dd75/Noubel%20on%20Collective%20Intelligence.pdf
- Oakley, B. (2010). *Evil Genes: Why Rome Fell, Hitler Rose, Enron Failed, and My Sister Stole My Mother's Boyfriend*, Prometheus Books.
- O'Connor, P. (2008). User-generated content and travel: A case study on Tripadvisor.com, *Information and Communication Technologies in Tourism 2008* pp. 47–58, Springer.
- Olguín, D. O. and Pentland, A. (2010). Assessing group performance from collective behavior, *Proc. of the CSCW Workshop on Collective Intelligence in Organizations: Toward a Research Agenda*, Citeseer, Savannah, GA, USA, pp. 1–5.
- Olguín, D. O. and Pentland, A. S. (2007). Sociometric badges: State of the art and future applications, *Doctoral colloquium presented at IEEE 11th International Symposium on Wearable Computers, Boston, MA*.
- Oliver, J. (1993). Decision graphs – an extension of decision trees, *Proceedings of the 4th International Workshop on Artificial Intelligence and Statistics*, pp. 343–350. Extended version available as TR173, Department of Computer Science, Monash University, Clayton, Australia.
- Oliver, J., Dowe, D. L. and Wallace, C. S. (1992). Inferring decision graphs using the Minimum Message Length principle, *Proceedings of the Australian Joint Conference on Artificial Intelligence*, pp. 361–367.
- Oliver, J. and Wallace, C. S. (1991). Inferring decision graphs, *Proceedings of the 12th International Joint Conference on Artificial Intelligence (IJCAI-91), workshop 8*.

- Ooi, J. N. and Dowe, D. L. (2005). Inferring phylogenetic graphs of natural languages using minimum message length, *11th Conference of the Spanish Association for Artificial Intelligence (CAEPIA 2005)*, Vol. 1, pp. 143–152.
- OpenAI (2016). Universe: a software platform for measuring and training AI, GitHub repository: <https://github.com/openai/universe>.
URL: <https://universe.openai.com>
- Oppy, G. and Dowe, D. L. (2011). The Turing Test, in E. N. Zalta (ed.), *Stanford Encyclopedia of Philosophy*, Stanford University.
URL: <http://plato.stanford.edu/entries/turing-test/>
- Orseau, L. and Ring, M. (2012). Space-time embedded intelligence, *International Conference on Artificial General Intelligence (AGI'12)*, Vol. 7716 of *Lecture Notes in Computer Science (LNCS)*, Springer, Oxford, UK, pp. 209–218.
- Ostrom, E. (1998). A behavioral approach to the rational choice theory of collective action: Presidential address, American Political Science Association, 1997, *American Political Science Review* **92**(1): 1–22, Cambridge University Press.
- Ostrom, E. (2014). Collective action and the evolution of social norms, *Journal of Natural Resources Policy Research* **6**(4): 235–252, Taylor & Francis.
- Pan, W., Dong, W., Cebrian, M., Kim, T., Fowler, J. H. and Pentland, A. S. (2012). Modeling dynamical influence in human interaction: Using data to make better inferences about influence within social systems, *Signal Processing Magazine, IEEE* **29**(2): 77–86, IEEE.
- Panait, L. and Luke, S. (2005). Cooperative multi-agent learning: The state of the art, *Autonomous Agents and Multi-Agent Systems* **11**(3): 387–434, Kluwer Academic Publishers.
URL: <http://dx.doi.org/10.1007/s10458-005-2631-2>
- Papentin, F. (1983). Binary sequences. I. complexity, *Information Sciences* **31**(1): 1–14, Elsevier.
- Parpinelli, R. S. and Lopes, H. S. (2011). New inspirations in swarm intelligence: a survey, *International Journal of Bio-Inspired Computation* **3**(1): 1–16, Inderscience Publishers Ltd.
- Parsons, S., Rodriguez-Aguilar, J. A. and Klein, M. (2011). Auctions and bidding: A guide for computer scientists, *ACM* **43**(2): 1–59, ACM.
URL: <http://doi.acm.org/10.1145/1883612.1883617>
- Passino, K. M. (2002). Biomimicry of bacterial foraging for distributed optimization and control, *Control Systems, IEEE* **22**(3): 52–67, IEEE.
- Pentland, A. (2006). Collective intelligence, *Computational Intelligence Magazine, IEEE* **1**(3): 9–12.

- Pentland, A. (2007). On the collective nature of human intelligence., *Adaptive Behaviour* **15**(2): 189–198.
URL: <http://dblp.uni-trier.de/db/journals/adb/adb15.html#Pentland07>
- Pfeifer, R. and Scheier, C. (2001). *Understanding intelligence*, MIT press.
- Pham, D. T., Ghanbarzadeh, A., Koc, E., Otri, S., Rahim, S. and Zaidi, M. (2006). The Bees Algorithm, A Novel Tool for Complex Optimisation Problems, *Proceedings of the 2nd International Virtual Conference on Intelligent Production Machines and Systems (IPROMS 2006)*, Elsevier, pp. 454–461.
- Piaget, J. (1952). *The origins of intelligence in children*, International Universities Press, New York. Translated by Margaret Cook.
- Plotnik, J. M., Lair, R., Suphachoksahakun, W. and De Waal, F. B. (2011). Elephants know when they need a helping trunk in a cooperative task, *Proceedings of the National Academy of Sciences* **108**(12): 5116–5121, National Academy of Sciences.
- Poesio, M., Chamberlain, J., Kruschwitz, U., Robaldo, L. and Ducceschi, L. (2015). Phrase detectives: Utilizing collective intelligence for internet-scale language resource creation (extended abstract), *Proceedings 24th IJCAI, Buenos Aires, Argentina*, pp. 4202–4206.
URL: <http://ijcai.org/papers15/Abstracts/IJCAI15-596.html>
- Popkin, G. (2016). The physics of life, *Nature* **529**(7584): 16–18, Nature Publishing Group.
- Potter, M. A. and De Jong, K. A. (2000). Cooperative coevolution: An architecture for evolving coadapted subcomponents, *Evolutionary Computation* **8**(1): 1–29, MIT Press.
- Poundstone, W. (2011). *Prisoner's dilemma*, Anchor.
- Prigogine, I. and Nicolis, G. (1989). *Exploring complexity: An introduction*, WH Freeman.
- Rabanal, P., Rodríguez, I. and Rubio, F. (2007). Using river formation dynamics to design heuristic algorithms, *International Conference on Unconventional Computation*, Springer, pp. 163–177.
- Raddick, M. J., Bracey, G., Gay, P. L., Lintott, C. J., Murray, P., Schawinski, K., Szalay, A. S. and Vandenberg, J. (2010). Galaxy Zoo: Exploring the motivations of citizen science volunteers, *Astronomy Education Review* **9**(1), The American Astronomical Society. 10.3847/AER2009036.
URL: <http://portico.org/stable?au=pgg3ztfdp8z>
- Ragni, M., Stahl, P. and Fangmeier, T. (2011). Cognitive complexity in matrix reasoning tasks, In B. Kokinov, A. Karmiloff-Smith, & N. J. Nersessian (Eds), *European Perspectives on Cognitive Science: Proceedings of the European conference on cognitive science (ECAI 2011)*, New Bulgarian University Press, Sofia.
- Rashedi, E., Nezamabadi-Pour, H. and Saryazdi, S. (2009). GSA: a gravitational search algorithm, *Information Sciences* **179**(13): 2232–2248, Elsevier.

- Raven, J. C. and Court, J. H. (1998). *Raven's progressive matrices and vocabulary scales*, Oxford Psychologists Press, Oxford, UK.
- Reid, C. R., Lutz, M. J., Powell, S., Kao, A. B., Couzin, I. D. and Garnier, S. (2015). Army ants dynamically adjust living bridges in response to a cost–benefit trade-off, *Proceedings of the National Academy of Sciences* **112**(49): 15113–15118, National Academy of Sciences.
- Reynolds, C. W. (1987). Flocks, herds and schools: A distributed behavioral model, *SIGGRAPH Computer Graphics* **21**(4): 25–34. (ACM SIGGRAPH '87 Conference Proceedings, Anaheim, California, July 1987).
- Reznikova, Z. I. (2007). *Animal intelligence: from individual to social cognition*, Cambridge University Press.
- Ring, M. and Orseau, L. (2011). Delusion, survival, and intelligent agents, *International Conference on Artificial General Intelligence*, Vol. 6830 of *Lecture Notes in Computer Science (LNCS)*, Springer, pp. 11–20.
- Rissanen, J. (1978). Modeling by shortest data description, *Automatica* **14**(5): 465–471, Elsevier.
- Roid, G. H. (2003). *Stanford-Binet intelligence scales*, Riverside Publishing Itasca, IL.
- Rong, Z., Yang, H.-X. and Wang, W.-X. (2010). Feedback reciprocity mechanism promotes the cooperation of highly clustered scale-free networks, *Physical Review E* **82**(4): 047101, APS.
- Rosa, M., Feyereisl, J. and The GoodAI team (2016). A framework for searching for general artificial intelligence, *arXiv preprint arXiv:1611.00685*.
URL: <https://arxiv.org/pdf/1611.00685.pdf>
- Roth, G. and Dicke, U. (2005). Evolution of the brain and intelligence, *Trends in Cognitive Sciences* **9**(5): 250–257, Elsevier.
- Rummery, G. A. and Niranjan, M. (1994). On-line Q-learning using connectionist systems, *Technical Report CUED/F-INFENG/TR 166*, University of Cambridge, Department of Engineering.
- Rust, J. and Golombok, S. (2014). *Modern psychometrics: The science of psychological assessment*, Routledge.
- Ryabko, B. and Reznikova, Z. (2009). The use of ideas of information theory for studying language and intelligence in ants, *Entropy* **11**(4): 836–853, Molecular Diversity Preservation International.
- Salminen, J. (2012). Collective intelligence in humans: A literature review, *arXiv preprint arXiv:1204.3401*.
URL: <https://arxiv.org/ftp/arxiv/papers/1204/1204.3401.pdf>

- Sanghi, P. and Dowe, D. L. (2003). A computer program capable of passing I.Q. tests, in P. Slezak (ed.), *Proceedings of the Joint International Conference on Cognitive Science, 4th ICCS International Conference on Cognitive Science & 7th ASCS Australasian Society for Cognitive Science (ICCS/ASCS-2003)*, Sydney, NSW, Australia, pp. 570–575.
- Satopää, V. A., Baron, J., Foster, D. P., Mellers, B. A., Tetlock, P. E. and Ungar, L. H. (2014). Combining multiple probability predictions using a simple logit model, *International Journal of Forecasting* **30**(2): 344–356, Elsevier.
- Schaie, K. W. (1979). Primary mental abilities in adulthood: an exploration in the development of psychometric intelligence, *Life-span Development and Behavior*. In P.B. Baltes & O. G. Brim Jr. (Eds) **3**: 67–115, New York: Academic Press.
- Schelling, T. C. (2006). *Micromotives and macrobehavior*, WW Norton & Company.
- Schmidt, D. F. (2008). *Minimum message length inference of autoregressive moving average models (PhD. thesis)*, Faculty of IT, Monash University.
- Schölkopf, B. (2015). Artificial intelligence: Learning to see and act, *Nature* **518**(7540): 486–487, Nature Research.
- Seeme, F. B. and Green, D. G. (2016). Pluralistic ignorance: Emergence and hypotheses testing in a multi-agent system, *Proc. of IEEE International Joint Conference on Neural Networks (IJCNN)*, pp. 5269–5274.
- Serugendo, G. D. M., Irit, M.-P. and Karageorgos, A. (2006). Self-organisation and emergence in MAS: An overview, *Informatica* **30**(1): 45–54.
- Service, E. T., Thurstone, L. L. and Thurstone, T. G. (1952). *American council on education psychological examination for college freshmen*, Cooperative Test Division, Educational Testing Service.
- Shah-Hosseini, H. (2009). The intelligent water drops algorithm: a nature-inspired swarm-based optimization algorithm, *International Journal of Bio-Inspired Computation* **1**(1-2): 71–79, Inderscience Publishers.
- Shannon, C. E. (1948). A mathematical theory of communication, *The Bell System Technical Journal* **27**(3): 379–423.
- Shapley, L. and Grofman, B. (1984). Optimizing group judgmental accuracy in the presence of interdependencies, *Public Choice* **43**(3): 329–343, Springer.
- Shettleworth, S. J. (2010). Clever animals and killjoy explanations in comparative psychology, *Trends in Cognitive Sciences* **14**(11): 477–481, Elsevier.
- Shultz, S. and Dunbar, R. (2006). Both social and ecological factors predict ungulate brain size, *Proceedings of the Royal Society of London B: Biological Sciences* **273**(1583): 207–215, The Royal Society.

- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M. et al. (2016). Mastering the game of Go with deep neural networks and tree search, *Nature* **529**(7587): 484–489, Nature Publishing Group.
- Simon, H. A. (1982). *Models of bounded rationality: Empirically grounded economic reason*, Vol. 3, MIT press.
- Simon, H. A. and Kotovsky, K. (1963). Human acquisition of concepts for sequential patterns., *Psychological Review* **70**(6): 534, American Psychological Association.
- Skyrms, B. and Pemantle, R. (2009). A dynamic model of social network formation, *Adaptive Networks: Theory, Models and Applications*, Springer, pp. 231–251. Chapter based on B. Skyrms and R. Pemantle, “A Dynamic Model of Social Network Formation” PNAS 97 (16), 9340-9346 (2000).
- Sloan, H. L., Good, M. and Dunnett, S. B. (2006). Double dissociation between hippocampal and prefrontal lesions on an operant delayed matching task and a water maze reference memory task, *Behavioural Brain Research* **171**(1): 116–126, Elsevier.
- Solomonoff, R. J. (1960). A preliminary report on a general theory of inductive inference (Report ZTB-138), *Cambridge, MA: Zator Co* **131**.
- Solomonoff, R. J. (1964a). A formal theory of inductive inference. part I, *Information and Control* **7**(1): 1–22, Elsevier.
- Solomonoff, R. J. (1964b). A formal theory of inductive inference. part II, *Information and Control* **7**(2): 224–254, Elsevier.
- Solomonoff, R. J. (1967). Inductive inference research: status, Spring 1967. RTB 154, Rockford Research, Inc., 140 1/2 Mt. Auburn St., Cambridge, Mass. 0213.
- Solomonoff, R. J. (1986). The application of algorithmic probability to problems in artificial intelligence, in L. Kanal and J. Lemmer (eds), *Uncertainty in Artificial Intelligence*, Elsevier Science Publishers, B.V., pp. 473–491. Also in: M. Kochen and H.M. Hastings, Eds., *Advances in Cognitive Science*, AAAS Selected Symposia Series, AAAS, Washington, D.C., pp. 210-227, 1988.
- Spearman, C. (1904). “General Intelligence”, objectively determined and measured, *The American Journal of Psychology* **15**(2): 201–292, JSTOR.
- Strannegård, C., Nizamani, A. R., Sjöberg, A. and Engström, F. (2013). Bounded Kolmogorov complexity based on cognitive models, *International Conference on Artificial General Intelligence*, Springer, pp. 130–139.
- Strogatz, S. H. (2001). Exploring complex networks, *Nature* **410**(6825): 268–276, Nature Publishing Group.
- Surowiecki, J. (2005). *The Wisdom of Crowds*, Anchor.

- Swain, D. T., Couzin, I. D. and Leonard, N. E. (2012). Real-time feedback-controlled robotic fish for behavioral experiments with fish schools, *Proceedings of the IEEE* **100**(1): 150–163, IEEE.
- Synnaeve, G., Nardelli, N., Auvolat, A., Chintala, S., Lacroix, T., Lin, Z., Richoux, F. and Usunier, N. (2016). Torchcraft: a library for machine learning research on real-time strategy games, *arXiv preprint arXiv:1611.00625*.
URL: <https://arxiv.org/pdf/1611.00625.pdf>
- Tan, P. J. and Dowe, D. L. (2002). MML inference of decision graphs with multi-way joins, *Proceedings of the 15th Australian Joint Conference on Artificial Intelligence*, Vol. 2557 of *Lecture Notes in Artificial Intelligence (LNAI)*, Springer, pp. 131–142.
- Tan, P. J. and Dowe, D. L. (2003). MML inference of decision graphs with multi-way joins and dynamic attributes, *Australasian Joint Conference on Artificial Intelligence*, Vol. 2903 of *Lecture Notes in Artificial Intelligence (LNAI)*, Springer, pp. 269–281.
- Tan, P. J. and Dowe, D. L. (2004). MML inference of oblique decision trees, *Australasian Joint Conference on Artificial Intelligence*, Vol. 3339 of *Lecture Notes in Artificial Intelligence (LNAI)*, Springer, pp. 1082–1088.
- Tanaka, K., Okamoto, K. and Tanaka, H. (2003). *Manual of the new Tanaka B Intelligence Scale*, Tokyo: Kaneko Shobo.
- Tapscott, D. and Williams, A. D. (2008). *Wikinomics: How mass collaboration changes everything*, Penguin.
- Thorndike, E. L. (1965). *Animal intelligence: Experimental studies*, Transaction Publishers.
- Thorndike, E. L., Bregman, E. O., Cobb, M. V. and Woodyard, E. (1926). *The measurement of intelligence.*, Teachers College Bureau of Publications.
- Thurstone, L. L. (1938). *Primary mental abilities*, Chicago Press.
- Thurstone, L. L. and Thurstone, T. G. (1941). *Factorial studies of intelligence: Psychometric Monographs*, Psychometric Society, issue 2, University of Chicago Press.
- Togelius, J., Shaker, N., Karakovskiy, S. and Yannakakis, G. N. (2013). The mario AI championship 2009-2012, *AI Magazine* **34**(3): 89–92.
- Tolman, E. C. (1948). Cognitive maps in rats and men, *Psychological Review* **55**(4): 189, American Psychological Association.
- Tomasello, M. and Call, J. (1997). *Primate cognition*, Oxford University Press, USA.
- Torsello, A. and Dowe, D. L. (2008a). Learning a generative model for structural representations, *Proceedings of the 21st Australasian Joint Conference on Artificial Intelligence*, Vol. 5360 of *Lecture Notes in Artificial Intelligence (LNAI)*, Springer, Auckland, NZ, pp. 573–583.

- Torsello, A. and Dowe, D. L. (2008b). Supervised learning of a generative model for edge-weighted graphs, *Proceedings of the 19th International Conference on Pattern Recognition, 2008 (ICPR 2008)*, IEEE Catalog Number: CFP08182, Tampa, Florida, U.S.A., pp. 1–4.
- Tran, Q. and Tian, Y. (2013). Organizational structure: Influencing factors and impact on a firm, *American Journal of Industrial and Business Management* **3**(2): 229–236, Scientific Research Publishing.
- Tumer, K. and Wolpert, D. (2004). A survey of collectives, *In collectives and the design of complex systems*, Springer, pp. 1–42.
- Turing, A. M. (1950). Computing machinery and intelligence, *Mind* **59**(236): 433–460, JSTOR.
- Uno, Y., Mizukami, H., Ando, M., Yukihiro, R., Iwasaki, Y. and Ozaki, N. (2014). Reliability and validity of the New Tanaka B Intelligence Scale scores: A group intelligence test, *PLoS One* **9**(6): 1–5, Public Library of Science.
URL: <http://dx.doi.org/10.1371/journal.pone.0100262>
- Vaarandi, R. (2004). A breadth-first algorithm for mining frequent patterns from event logs, *Intelligence in Communication Systems*, Springer, pp. 293–308.
- Valiant, L. G. (1984). A theory of the learnable, *Communications of the ACM* **27**(11): 1134–1142, ACM.
- Van Veldhuizen, D. A. and Lamont, G. B. (2000). Multiobjective evolutionary algorithms: Analyzing the state-of-the-art, *Evolutionary Computation* **8**(2): 125–147, MIT Press.
- Visser, G., Dale, P., Dowe, D. L., Ndoen, E., Dale, M. and Sipe, N. (2012). A novel approach for modeling malaria incidence using complex categorical household data: The Minimum Message Length (MML) method applied to Indonesian data, *Computational Ecology and Software* **2**(3): 140–159, International Academy of Ecology and Environmental Sciences (IAEES).
- Visser, G. and Dowe, D. L. (2007). Minimum Message Length clustering of spatially-correlated data with varying inter-class penalties, *Proceedings of the 6th IEEE/ACIS International Conference on Computer and Information Science (ICIS 2007)*, IEEE, pp. 17–22.
- Viswanathan, M., Wallace, C. S., Dowe, D. L. and Korb, K. B. (1999). Finding outpoints in noisy binary sequences – a revised empirical evaluation, *Australasian Joint Conference on Artificial Intelligence*, Vol. 1747 of *Lecture Notes in Artificial Intelligence (LNAI)*, Springer, Sydney, Australia, pp. 405–416.
- Von Ahn, L., Maurer, B., McMillen, C., Abraham, D. and Blum, M. (2008). reCAPTCHA: Human-based character recognition via web security measures, *Science* **321**(5895): 1465–1468, American Association for the Advancement of Science.

- Waldrop, M. M. (1993). *Complexity: The emerging science at the edge of order and chaos*, Simon and Schuster.
- Wallace, C. S. (1986). An improved program for classification, *Proceedings of the 9th Australian Computer Science Conference (ACSC-9)*, Vol. 8, pp. 357–366.
- Wallace, C. S. (1998). Intrinsic classification of spatially correlated data, *The Computer Journal* **41**(8): 602–611, British Computer Society.
- Wallace, C. S. (2005). *Statistical and Inductive Inference by Minimum Message Length*, Springer-Verlag.
- Wallace, C. S. and Boulton, D. M. (1968). An information measure for classification, *The Computer Journal* **11**(2): 185–194, British Computer Society.
- Wallace, C. S. and Boulton, D. M. (1975). An invariant Bayes method for point estimation, *Classification Society Bulletin* **3**(3): 11–34.
- Wallace, C. S. and Dowe, D. L. (1993). MML estimation of the von Mises concentration parameter, *Technical Report Technical report 93/193*, Department of Computer Science, Monash University, Melbourne, Australia.
- Wallace, C. S. and Dowe, D. L. (1994a). Estimation of the von Mises concentration parameter using Minimum Message Length, *Proc. 12th Australian Statistical Society Conference (Abstract)*, Monash University, Australia.
- Wallace, C. S. and Dowe, D. L. (1994b). Intrinsic classification by MML – the Snob program, in C. Zhang, J. Debenham and D. Lukose (eds), *Proceedings of the 7th Australian Joint Conference on Artificial Intelligence: sowing the seeds for the future*, World Scientific, Armidale, Australia, pp. 37–44.
- Wallace, C. S. and Dowe, D. L. (1999). Minimum message length and Kolmogorov complexity, *The Computer Journal* **42**(4): 270–283, British Computer Society. Special issue on Kolmogorov complexity.
- Wallace, C. S. and Dowe, D. L. (2000). MML clustering of multi-state, Poisson, von Mises circular and Gaussian distributions, *Statistics and Computing* **10**(1): 73–83, Springer.
- Wallace, C. S. and Patrick, J. D. (1993). Coding decision trees, *Machine Learning* **11**(1): 7–22, Springer.
- Wasserman, E. A. and Zentall, T. R. (2006). *Comparative cognition: Experimental explorations of animal intelligence*, Oxford University Press, USA.
- Watkins, C. J. C. H. and Dayan, P. (1992). Technical note: Q-learning, *Machine Learning* **8**(3-4): 279–292, Kluwer Academic Publishers.
URL: <http://dx.doi.org/10.1007/BF00992698>
- Watts, D. J. (2004). *Six degrees: The science of a connected age*, WW Norton & Company.

- Watts, D. J. and Strogatz, S. H. (1998). Collective dynamics of “small-world” networks, *Nature* **393**(6684): 440–442, Nature Publishing Group.
- Wechsler, D. (2008). *Wechsler adult intelligence scale-fourth*, San Antonio: Pearson.
- Wechsler, D. and Hardesty, A. (1964). *Die Messung der Intelligenz Erwachsener*, Huber.
- Wedde, H. F., Farooq, M. and Zhang, Y. (2004). Beehive: An efficient fault-tolerant routing algorithm inspired by honey bee behavior, in M. Dorigo, M. Birattari, C. Blum, L. M. Gambardella, F. Mondada and T. Stützle (eds), *4th International Workshop on Ant Colony Optimization and Swarm Intelligence, ANTS 2004*, Springer, Brussels, Belgium, pp. 83–94.
- Weisbuch, G., Kirman, A. and Herreiner, D. (2000). Market organisation and trading relationships, *The Economic Journal* **110**(463): 411–436, Wiley Online Library.
- Weisstein, E. W. (2015). Moore neighborhood, from mathworld - a Wolfram web resource. Last accessed: 13/03/2017.
URL: <http://mathworld.wolfram.com/MooreNeighborhood.html>
- Weschler, D. (1971). Concept of collective intelligence., *American Psychologist* **26**(10): 904, American Psychological Association.
- Weyns, D., Steegmans, E. and Holvoet, T. (2004). Towards active perception in situated multiagent systems, *Applied Artificial Intelligence* **18**(9-10): 867–883.
URL: <http://dx.doi.org/10.1080/08839510490509063>
- Wilensky, U. (1999). NetLogo: Center for connected learning and computer-based modeling, Northwestern University, Evanston, IL.
URL: <http://ccl.northwestern.edu/netlogo/>
- Wolfers, J. and Zitzewitz, E. (2004). Prediction markets, *The Journal of Economic Perspectives* **18**(2): 107–126, American Economic Association.
- Wolpert, D. (2004). Theory of collective intelligence, *Collectives and the design of complex systems*, Springer, pp. 43–106.
- Wolpert, D. H. and Macready, W. G. (1997). No free lunch theorems for optimization, *IEEE Transactions on Evolutionary Computation* **1**(1): 67–82, IEEE.
- Wolpert, D. H. and Tumer, K. (1999). An introduction to collective intelligence, *Technical Report NASA-ARC-IC-99-63*. arXiv preprint cs/9908014.
URL: <http://arxiv.org/abs/cs.LG/9908014>
- Wolpert, D., Tumer, K., Server and Nasa Technical Reports Server (Ntrs) (2013). *A Survey of Collective Intelligence*, reprint edn, BiblioLife.
URL: <http://books.google.com.au/books?id=TMdOngEACAAJ>

- Woodley, M. A. and Bell, E. (2011). Is collective intelligence (mostly) the general factor of personality? A comment on Woolley, Chabris, Pentland, Hashmi and Malone (2010), *Intelligence* **39**(2-3): 79–81.
URL: <http://www.sciencedirect.com/science/article/pii/S0160289611000201>
- Wooldridge, M. and Jennings, N. R. (1995). Intelligent agents: Theory and practice, *The Knowledge Engineering Review* **10**(2): 115–152, Cambridge University Press.
- Woolley, A. W., Chabris, C. F., Pentland, A., Hashmi, N. and Malone, T. W. (2010). Evidence for a collective intelligence factor in the performance of human groups, *Science* **330**(6004): 686–688, American Association for the Advancement of Science.
URL: <http://www.sciencemag.org/content/330/6004/686.abstract>
- Khafa, F. and Bessis, N. (2014). *Inter-cooperative Collective Intelligence: Techniques and Applications*, Springer.
- Yang, X.-S. (2009). Firefly algorithms for multimodal optimization, *International Symposium on Stochastic algorithms*, Springer, pp. 169–178.
- Yang, X.-S. (2010a). *Nature-inspired metaheuristic algorithms*, Luniver Press.
- Yang, X.-S. (2010b). A new metaheuristic bat-inspired algorithm, in J. R. Gonzalez (ed.), *Nature inspired cooperative strategies for optimization (NICSO 2010)*, Vol. 284 of *Studies in Computational Intelligence*, Springer, pp. 65–74.
- Yang, X.-S. and Deb, S. (2009). Cuckoo search via Lévy flights, *World Congress on Nature & Biologically Inspired Computing, NaBIC 2009*, IEEE, pp. 210–214.
- Yu, C.-H., Werfel, J. and Nagpal, R. (2010). Collective decision-making in multi-agent systems by implicit leadership, in W. van der Hoek, G. A. Kaminka, Y. Lespérance, M. Luck and S. Sen (eds), *Proceedings of the 9th International Conference on Autonomous Agents and Multiagent Systems*, Vol. 3 of *AAMAS'10*, International Foundation for Autonomous Agents and Multiagent Systems, Toronto, Canada, pp. 1189–1196.
URL: <https://dash.harvard.edu/handle/1/9943235>
- Zhang, S., Mizutani, A. and Srinivasan, M. V. (2000). Maze navigation by honeybees: learning path regularity, *Learning & Memory* **7**(6): 363–374, Cold Spring Harbor Lab.
- Zimmermann, M. G., Eguiluz, V. M. and San Miguel, M. (2001). Cooperation, adaptation and the emergence of leadership, *Economics with heterogeneous interacting agents*, Springer, pp. 73–86.
- Ziv, J. and Lempel, A. (1978). Compression of individual sequences via variable-rate coding, *IEEE transactions on Information Theory* **24**(5): 530–536, IEEE.