



GLOBAL EDITION



Stuart  
**Russell**  
Peter  
**Norvig**

# Artificial Intelligence

## A Modern Approach

*Fourth Edition*



# Artificial Intelligence

A Modern Approach

*Fourth Edition*

*Global Edition*



**PEARSON SERIES  
IN ARTIFICIAL INTELLIGENCE**

*Stuart Russell and Peter Norvig, Editors*

FORSYTH & PONCE

GRAHAM

JURAFSKY & MARTIN

NEAPOLITAN

RUSSELL & NORVIG

*Computer Vision: A Modern Approach, 2nd ed.*

*ANSI Common Lisp*

*Speech and Language Processing, 2nd ed.*

*Learning Bayesian Networks*

*Artificial Intelligence: A Modern Approach, 4th ed.*

# Artificial Intelligence

A Modern Approach

*Fourth Edition*

*Global Edition*

Stuart J. Russell and Peter Norvig

*Contributing writers:*

Ming-Wei Chang

Jacob Devlin

Anca Dragan

David Forsyth

Ian Goodfellow

Jitendra M. Malik

Vikash Mansinghka

Judea Pearl

Michael Wooldridge



Cover Image credits: Alan Turing: Science History Images/Alamy Stock Photo; Statue of Aristotle: Panos Karas/Shutterstock; Ada Lovelace – Pictorial Press Ltd/Alamy Stock Photo; Autonomous cars: Andrey Suslov/Shutterstock; Atlas Robot: Boston Dynamics, Inc.; Berkeley Campanile and Golden Gate Bridge: Ben Chu/Shutterstock; Background ghosted nodes: Eugene Sergeev/Alamy Stock Photo; Chess board with chess figure: Titania/Shutterstock; Mars Rover: Stocktrek Images, Inc./Alamy Stock Photo; Kasparov: KATHY WILLENS/AP Images

*Pearson Education Limited*

KAO Two  
KAO Park  
Hockham Way  
Harlow  
CM17 9SR  
United Kingdom

and Associated Companies throughout the world

*Visit us on the World Wide Web at:* [www.pearsonglobaleditions.com](http://www.pearsonglobaleditions.com)

Please contact <https://support.pearson.com/getsupport/s/contactsupport> with any queries on this content

© Pearson Education Limited 2022

The rights of Stuart Russell and Peter Norvig to be identified as the authors of this work have been asserted by them in accordance with the Copyright, Designs and Patents Act 1988.

*Authorized adaptation from the United States edition, entitled Artificial Intelligence: A Modern Approach, 4th Edition, ISBN 978-0-13-461099-3 by Stuart J. Russell and Peter Norvig, published by Pearson Education © 2021.*

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, without either the prior written permission of the publisher or a license permitting restricted copying in the United Kingdom issued by the Copyright Licensing Agency Ltd, Saffron House, 6–10 Kirby Street, London EC1N 8TS. For information regarding permissions, request forms and the appropriate contacts within the Pearson Education Global Rights & Permissions department, please visit [www.pearsoned.com/permissions/](http://www.pearsoned.com/permissions/).

All trademarks used herein are the property of their respective owners. The use of any trademark in this text does not vest in the author or publisher any trademark ownership rights in such trademarks, nor does the use of such trademarks imply any affiliation with or endorsement of this book by such owners.

PEARSON, ALWAYS LEARNING, and MYLAB are exclusive trademarks in the U.S. and/or other countries owned by Pearson Education, Inc. or its affiliates.

Unless otherwise indicated herein, any third-party trademarks that may appear in this work are the property of their respective owners and any references to third-party trademarks, logos or other trade dress are for demonstrative or descriptive purposes only. Such references are not intended to imply any sponsorship, endorsement, authorization, or promotion of Pearson's products by the owners of such marks, or any relationship between the owner and Pearson Education, Inc. or its affiliates, authors, licensees or distributors.

**ISBN 10:** 1-292-40113-3

**ISBN 13:** 978-1-292-40113-3

**eBook ISBN 13:** 978-1-292-40117-1

#### **British Library Cataloguing-in-Publication Data**

A catalogue record for this book is available from the British Library

Typeset by SPi Global

eBook formatted by B2R Technologies Pvt. Ltd.

*For Loy, Gordon, Lucy, George, and Isaac — S.J.R.*

*For Kris, Isabella, and Juliet — P.N.*

*This page is intentionally left blank*

# Preface

**Artificial Intelligence** (AI) is a big field, and this is a big book. We have tried to explore the full breadth of the field, which encompasses logic, probability, and continuous mathematics; perception, reasoning, learning, and action; fairness, trust, social good, and safety; and applications that range from microelectronic devices to robotic planetary explorers to online services with billions of users.

The subtitle of this book is “A Modern Approach.” That means we have chosen to tell the story from a current perspective. We synthesize what is now known into a common framework, recasting early work using the ideas and terminology that are prevalent today. We apologize to those whose subfields are, as a result, less recognizable.

## New to this edition

This edition reflects the changes in AI since the last edition in 2010:

- We focus more on machine learning rather than hand-crafted knowledge engineering, due to the increased availability of data, computing resources, and new algorithms.
- Deep learning, probabilistic programming, and multiagent systems receive expanded coverage, each with their own chapter.
- The coverage of natural language understanding, robotics, and computer vision has been revised to reflect the impact of deep learning.
- The robotics chapter now includes robots that interact with humans and the application of reinforcement learning to robotics.
- Previously we defined the goal of AI as creating systems that try to maximize expected utility, where the specific utility information—the objective—is supplied by the human designers of the system. Now we no longer assume that the objective is fixed and known by the AI system; instead, the system may be uncertain about the true objectives of the humans on whose behalf it operates. It must learn what to maximize and must function appropriately even while uncertain about the objective.
- We increase coverage of the impact of AI on society, including the vital issues of ethics, fairness, trust, and safety.
- We have moved the exercises from the end of each chapter to an online site. This allows us to continuously add to, update, and improve the exercises, to meet the needs of instructors and to reflect advances in the field and in AI-related software tools.
- Overall, about 25% of the material in the book is brand new. The remaining 75% has been largely rewritten to present a more unified picture of the field. 22% of the citations in this edition are to works published after 2010.

## Overview of the book

The main unifying theme is the idea of an **intelligent agent**. We define AI as the study of agents that receive percepts from the environment and perform actions. Each such agent implements a function that maps percept sequences to actions, and we cover different ways to represent these functions, such as reactive agents, real-time planners, decision-theoretic

systems, and deep learning systems. We emphasize learning both as a construction method for competent systems and as a way of extending the reach of the designer into unknown environments. We treat robotics and vision not as independently defined problems, but as occurring in the service of achieving goals. We stress the importance of the task environment in determining the appropriate agent design.

Our primary aim is to convey the *ideas* that have emerged over the past seventy years of AI research and the past two millennia of related work. We have tried to avoid excessive formality in the presentation of these ideas, while retaining precision. We have included mathematical formulas and pseudocode algorithms to make the key ideas concrete; mathematical concepts and notation are described in Appendix A and our pseudocode is described in Appendix B.

This book is primarily intended for use in an undergraduate course or course sequence. The book has 29 chapters, each requiring about a week's worth of lectures, so working through the whole book requires a two-semester sequence. A one-semester course can use selected chapters to suit the interests of the instructor and students. The book can also be used in a graduate-level course (perhaps with the addition of some of the primary sources suggested in the bibliographical notes), or for self-study or as a reference.

Throughout the book, *important points* are marked with a triangle icon in the margin. Wherever a new **term** is defined, it is also noted in the margin. Subsequent significant uses of the **term** are in bold, but not in the margin. We have included a comprehensive index and an extensive bibliography.

The only prerequisite is familiarity with basic concepts of computer science (algorithms, data structures, complexity) at a sophomore level. Freshman calculus and linear algebra are useful for some of the topics.

## Online resources

Online resources are available through [pearsonglobaleditions.com](http://pearsonglobaleditions.com). There you will find:

- Exercises, programming projects, and research projects. These are no longer at the end of each chapter; they are online only. Within the book, we refer to an online exercise with a name like “Exercise 6.NARY.” Instructions on the Web site allow you to find exercises by name or by topic.
- Implementations of the algorithms in the book in Python, Java, and other programming languages.
- Supplementary material and links for students and instructors.
- Instructions on how to report errors in the book in the likely event that some exist.

## Book cover

The cover depicts the final position from the decisive game 6 of the 1997 chess match in which the program Deep Blue defeated Garry Kasparov (playing Black), making this the first time a computer had beaten a world champion in a chess match. Kasparov is shown at the top. To his right is a pivotal position from the second game of the historic Go match between former world champion Lee Sedol and DeepMind's ALPHAGO program. Move 37 by ALPHAGO violated centuries of Go orthodoxy and was immediately seen by human experts

as an embarrassing mistake, but it turned out to be a winning move. At top left is an Atlas humanoid robot built by Boston Dynamics. A depiction of a self-driving car sensing its environment appears between Ada Lovelace, the world's first computer programmer, and Alan Turing, whose fundamental work defined artificial intelligence. At the bottom of the chess board are a Mars Exploration Rover robot and a statue of Aristotle, who pioneered the study of logic; his planning algorithm from *De Motu Animalium* appears behind the authors' names. Behind the chess board is a probabilistic programming model used by the UN Comprehensive Nuclear-Test-Ban Treaty Organization for detecting nuclear explosions from seismic signals.

## Acknowledgments

It takes a global village to make a book. Over 600 people read parts of the book and made suggestions for improvement. The complete list is at [pearsonglobaleditions.com](http://pearsonglobaleditions.com); we are grateful to all of them. We have space here to mention only a few especially important contributors. First the contributing writers:

- Judea Pearl (Section 13.5, Causal Networks);
- Michael Wooldridge (Chapter 17, Multiagent Decision Making);
- Vikash Mansinghka (Section 18.4, Programs as Probability Models);
- Ian Goodfellow (Chapter 22, Deep Learning);
- Jacob Devlin and Mei-Wing Chang (Chapter 25, Deep Learning for Natural Language Processing);
- Anca Dragan (Chapter 26, Robotics);
- Jitendra Malik and David Forsyth (Chapter 27, Computer Vision).

Then some key roles:

- Cynthia Yeung and Malika Cantor (project management);
- Julie Sussman and Tom Galloway (copyediting and writing suggestions);
- Omari Stephens (illustrations);
- Tracy Johnson (editor);
- Erin Ault and Rose Kernan (cover and color conversion);
- Nalin Chhibber, Sam Goto, Raymond de Lacaze, Ravi Mohan, Ciaran O'Reilly, Amit Patel, Dragomir Radiv, and Samagra Sharma (online code development and mentoring);
- Google Summer of Code students (online code development).

**Stuart would like to thank** his wife, Loy Sheflott, for her endless patience and boundless wisdom. He hopes that Gordon, Lucy, George, and Isaac will soon be reading this book after they have forgiven him for working so long on it. RUGS (Russell's Unusual Group of Students) have been unusually helpful, as always.

**Peter would like to thank** his parents (Torsten and Gerda) for getting him started, and his wife (Kris), children (Bella and Juliet), colleagues, boss, and friends for encouraging and tolerating him through the long hours of writing and rewriting.

## About the Authors

**Stuart Russell** was born in 1962 in Portsmouth, England. He received his B.A. with first-class honours in physics from Oxford University in 1982, and his Ph.D. in computer science from Stanford in 1986. He then joined the faculty of the University of California at Berkeley, where he is a professor and former chair of computer science, director of the Center for Human-Compatible AI, and holder of the Smith–Zadeh Chair in Engineering. In 1990, he received the Presidential Young Investigator Award of the National Science Foundation, and in 1995 he was cowinner of the Computers and Thought Award. He is a Fellow of the American Association for Artificial Intelligence, the Association for Computing Machinery, and the American Association for the Advancement of Science, an Honorary Fellow of Wadham College, Oxford, and an Andrew Carnegie Fellow. He held the Chaire Blaise Pascal in Paris from 2012 to 2014. He has published over 300 papers on a wide range of topics in artificial intelligence. His other books include *The Use of Knowledge in Analogy and Induction*, *Do the Right Thing: Studies in Limited Rationality* (with Eric Wefald), and *Human Compatible: Artificial Intelligence and the Problem of Control*.

**Peter Norvig** is currently a Director of Research at Google, Inc., and was previously the director responsible for the core Web search algorithms. He co-taught an online AI class that signed up 160,000 students, helping to kick off the current round of massive open online classes. He was head of the Computational Sciences Division at NASA Ames Research Center, overseeing research and development in artificial intelligence and robotics. He received a B.S. in applied mathematics from Brown University and a Ph.D. in computer science from Berkeley. He has been a professor at the University of Southern California and a faculty member at Berkeley and Stanford. He is a Fellow of the American Association for Artificial Intelligence, the Association for Computing Machinery, the American Academy of Arts and Sciences, and the California Academy of Science. His other books are *Paradigms of AI Programming: Case Studies in Common Lisp*, *Verbomobil: A Translation System for Face-to-Face Dialog*, and *Intelligent Help Systems for UNIX*.

The two authors shared the inaugural AAAI/EAAI Outstanding Educator award in 2016.

# Contents

## I Artificial Intelligence

<b>1 Introduction</b>	<b>19</b>
1.1 What Is AI? . . . . .	19
1.2 The Foundations of Artificial Intelligence . . . . .	23
1.3 The History of Artificial Intelligence . . . . .	35
1.4 The State of the Art . . . . .	45
1.5 Risks and Benefits of AI . . . . .	49
Summary . . . . .	52
Bibliographical and Historical Notes . . . . .	53
<b>2 Intelligent Agents</b>	<b>54</b>
2.1 Agents and Environments . . . . .	54
2.2 Good Behavior: The Concept of Rationality . . . . .	57
2.3 The Nature of Environments . . . . .	60
2.4 The Structure of Agents . . . . .	65
Summary . . . . .	78
Bibliographical and Historical Notes . . . . .	78

## II Problem-solving

<b>3 Solving Problems by Searching</b>	<b>81</b>
3.1 Problem-Solving Agents . . . . .	81
3.2 Example Problems . . . . .	84
3.3 Search Algorithms . . . . .	89
3.4 Uninformed Search Strategies . . . . .	94
3.5 Informed (Heuristic) Search Strategies . . . . .	102
3.6 Heuristic Functions . . . . .	115
Summary . . . . .	122
Bibliographical and Historical Notes . . . . .	124
<b>4 Search in Complex Environments</b>	<b>128</b>
4.1 Local Search and Optimization Problems . . . . .	128
4.2 Local Search in Continuous Spaces . . . . .	137
4.3 Search with Nondeterministic Actions . . . . .	140
4.4 Search in Partially Observable Environments . . . . .	144
4.5 Online Search Agents and Unknown Environments . . . . .	152
Summary . . . . .	159
Bibliographical and Historical Notes . . . . .	160
<b>5 Constraint Satisfaction Problems</b>	<b>164</b>
5.1 Defining Constraint Satisfaction Problems . . . . .	164
5.2 Constraint Propagation: Inference in CSPs . . . . .	169

5.3	Backtracking Search for CSPs . . . . .	175
5.4	Local Search for CSPs . . . . .	181
5.5	The Structure of Problems . . . . .	183
	Summary . . . . .	187
	Bibliographical and Historical Notes . . . . .	188
<b>6</b>	<b>Adversarial Search and Games</b>	<b>192</b>
6.1	Game Theory . . . . .	192
6.2	Optimal Decisions in Games . . . . .	194
6.3	Heuristic Alpha–Beta Tree Search . . . . .	202
6.4	Monte Carlo Tree Search . . . . .	207
6.5	Stochastic Games . . . . .	210
6.6	Partially Observable Games . . . . .	214
6.7	Limitations of Game Search Algorithms . . . . .	219
	Summary . . . . .	220
	Bibliographical and Historical Notes . . . . .	221
 <b>III Knowledge, reasoning, and planning</b>		
<b>7</b>	<b>Logical Agents</b>	<b>226</b>
7.1	Knowledge-Based Agents . . . . .	227
7.2	The Wumpus World . . . . .	228
7.3	Logic . . . . .	232
7.4	Propositional Logic: A Very Simple Logic . . . . .	235
7.5	Propositional Theorem Proving . . . . .	240
7.6	Effective Propositional Model Checking . . . . .	250
7.7	Agents Based on Propositional Logic . . . . .	255
	Summary . . . . .	264
	Bibliographical and Historical Notes . . . . .	265
<b>8</b>	<b>First-Order Logic</b>	<b>269</b>
8.1	Representation Revisited . . . . .	269
8.2	Syntax and Semantics of First-Order Logic . . . . .	274
8.3	Using First-Order Logic . . . . .	283
8.4	Knowledge Engineering in First-Order Logic . . . . .	289
	Summary . . . . .	295
	Bibliographical and Historical Notes . . . . .	296
<b>9</b>	<b>Inference in First-Order Logic</b>	<b>298</b>
9.1	Propositional vs. First-Order Inference . . . . .	298
9.2	Unification and First-Order Inference . . . . .	300
9.3	Forward Chaining . . . . .	304
9.4	Backward Chaining . . . . .	311
9.5	Resolution . . . . .	316
	Summary . . . . .	327
	Bibliographical and Historical Notes . . . . .	328

<b>10 Knowledge Representation</b>	<b>332</b>
10.1 Ontological Engineering . . . . .	332
10.2 Categories and Objects . . . . .	335
10.3 Events . . . . .	340
10.4 Mental Objects and Modal Logic . . . . .	344
10.5 Reasoning Systems for Categories . . . . .	347
10.6 Reasoning with Default Information . . . . .	351
Summary . . . . .	355
Bibliographical and Historical Notes . . . . .	356
<b>11 Automated Planning</b>	<b>362</b>
11.1 Definition of Classical Planning . . . . .	362
11.2 Algorithms for Classical Planning . . . . .	366
11.3 Heuristics for Planning . . . . .	371
11.4 Hierarchical Planning . . . . .	374
11.5 Planning and Acting in Nondeterministic Domains . . . . .	383
11.6 Time, Schedules, and Resources . . . . .	392
11.7 Analysis of Planning Approaches . . . . .	396
Summary . . . . .	397
Bibliographical and Historical Notes . . . . .	398
<b>IV Uncertain knowledge and reasoning</b>	
<b>12 Quantifying Uncertainty</b>	<b>403</b>
12.1 Acting under Uncertainty . . . . .	403
12.2 Basic Probability Notation . . . . .	406
12.3 Inference Using Full Joint Distributions . . . . .	413
12.4 Independence . . . . .	415
12.5 Bayes' Rule and Its Use . . . . .	417
12.6 Naive Bayes Models . . . . .	420
12.7 The Wumpus World Revisited . . . . .	422
Summary . . . . .	425
Bibliographical and Historical Notes . . . . .	426
<b>13 Probabilistic Reasoning</b>	<b>430</b>
13.1 Representing Knowledge in an Uncertain Domain . . . . .	430
13.2 The Semantics of Bayesian Networks . . . . .	432
13.3 Exact Inference in Bayesian Networks . . . . .	445
13.4 Approximate Inference for Bayesian Networks . . . . .	453
13.5 Causal Networks . . . . .	467
Summary . . . . .	471
Bibliographical and Historical Notes . . . . .	472
<b>14 Probabilistic Reasoning over Time</b>	<b>479</b>
14.1 Time and Uncertainty . . . . .	479
14.2 Inference in Temporal Models . . . . .	483

14.3	Hidden Markov Models . . . . .	491
14.4	Kalman Filters . . . . .	497
14.5	Dynamic Bayesian Networks . . . . .	503
	Summary . . . . .	514
	Bibliographical and Historical Notes . . . . .	515
<b>15</b>	<b>Making Simple Decisions</b>	<b>518</b>
15.1	Combining Beliefs and Desires under Uncertainty . . . . .	518
15.2	The Basis of Utility Theory . . . . .	519
15.3	Utility Functions . . . . .	522
15.4	Multiattribute Utility Functions . . . . .	530
15.5	Decision Networks . . . . .	534
15.6	The Value of Information . . . . .	537
15.7	Unknown Preferences . . . . .	543
	Summary . . . . .	547
	Bibliographical and Historical Notes . . . . .	547
<b>16</b>	<b>Making Complex Decisions</b>	<b>552</b>
16.1	Sequential Decision Problems . . . . .	552
16.2	Algorithms for MDPs . . . . .	562
16.3	Bandit Problems . . . . .	571
16.4	Partially Observable MDPs . . . . .	578
16.5	Algorithms for Solving POMDPs . . . . .	580
	Summary . . . . .	585
	Bibliographical and Historical Notes . . . . .	586
<b>17</b>	<b>Multiagent Decision Making</b>	<b>589</b>
17.1	Properties of Multiagent Environments . . . . .	589
17.2	Non-Cooperative Game Theory . . . . .	595
17.3	Cooperative Game Theory . . . . .	616
17.4	Making Collective Decisions . . . . .	622
	Summary . . . . .	635
	Bibliographical and Historical Notes . . . . .	636
<b>18</b>	<b>Probabilistic Programming</b>	<b>641</b>
18.1	Relational Probability Models . . . . .	642
18.2	Open-Universe Probability Models . . . . .	648
18.3	Keeping Track of a Complex World . . . . .	655
18.4	Programs as Probability Models . . . . .	660
	Summary . . . . .	664
	Bibliographical and Historical Notes . . . . .	665
<b>V</b>	<b>Machine Learning</b>	
<b>19</b>	<b>Learning from Examples</b>	<b>669</b>
19.1	Forms of Learning . . . . .	669

19.2	Supervised Learning . . . . .	671
19.3	Learning Decision Trees . . . . .	675
19.4	Model Selection and Optimization . . . . .	683
19.5	The Theory of Learning . . . . .	690
19.6	Linear Regression and Classification . . . . .	694
19.7	Nonparametric Models . . . . .	704
19.8	Ensemble Learning . . . . .	714
19.9	Developing Machine Learning Systems . . . . .	722
	Summary . . . . .	732
	Bibliographical and Historical Notes . . . . .	733
<b>20</b>	<b>Knowledge in Learning</b>	<b>739</b>
20.1	A Logical Formulation of Learning . . . . .	739
20.2	Knowledge in Learning . . . . .	747
20.3	Explanation-Based Learning . . . . .	750
20.4	Learning Using Relevance Information . . . . .	754
20.5	Inductive Logic Programming . . . . .	758
	Summary . . . . .	767
	Bibliographical and Historical Notes . . . . .	768
<b>21</b>	<b>Learning Probabilistic Models</b>	<b>772</b>
21.1	Statistical Learning . . . . .	772
21.2	Learning with Complete Data . . . . .	775
21.3	Learning with Hidden Variables: The EM Algorithm . . . . .	788
	Summary . . . . .	797
	Bibliographical and Historical Notes . . . . .	798
<b>22</b>	<b>Deep Learning</b>	<b>801</b>
22.1	Simple Feedforward Networks . . . . .	802
22.2	Computation Graphs for Deep Learning . . . . .	807
22.3	Convolutional Networks . . . . .	811
22.4	Learning Algorithms . . . . .	816
22.5	Generalization . . . . .	819
22.6	Recurrent Neural Networks . . . . .	823
22.7	Unsupervised Learning and Transfer Learning . . . . .	826
22.8	Applications . . . . .	833
	Summary . . . . .	835
	Bibliographical and Historical Notes . . . . .	836
<b>23</b>	<b>Reinforcement Learning</b>	<b>840</b>
23.1	Learning from Rewards . . . . .	840
23.2	Passive Reinforcement Learning . . . . .	842
23.3	Active Reinforcement Learning . . . . .	848
23.4	Generalization in Reinforcement Learning . . . . .	854
23.5	Policy Search . . . . .	861
23.6	Apprenticeship and Inverse Reinforcement Learning . . . . .	863

23.7 Applications of Reinforcement Learning . . . . .	866
Summary . . . . .	869
Bibliographical and Historical Notes . . . . .	870
<b>VI Communicating, perceiving, and acting</b>	
<b>24 Natural Language Processing</b>	<b>874</b>
24.1 Language Models . . . . .	874
24.2 Grammar . . . . .	884
24.3 Parsing . . . . .	886
24.4 Augmented Grammars . . . . .	892
24.5 Complications of Real Natural Language . . . . .	896
24.6 Natural Language Tasks . . . . .	900
Summary . . . . .	901
Bibliographical and Historical Notes . . . . .	902
<b>25 Deep Learning for Natural Language Processing</b>	<b>907</b>
25.1 Word Embeddings . . . . .	907
25.2 Recurrent Neural Networks for NLP . . . . .	911
25.3 Sequence-to-Sequence Models . . . . .	915
25.4 The Transformer Architecture . . . . .	919
25.5 Pretraining and Transfer Learning . . . . .	922
25.6 State of the art . . . . .	926
Summary . . . . .	929
Bibliographical and Historical Notes . . . . .	929
<b>26 Robotics</b>	<b>932</b>
26.1 Robots . . . . .	932
26.2 Robot Hardware . . . . .	933
26.3 What kind of problem is robotics solving? . . . . .	937
26.4 Robotic Perception . . . . .	938
26.5 Planning and Control . . . . .	945
26.6 Planning Uncertain Movements . . . . .	963
26.7 Reinforcement Learning in Robotics . . . . .	965
26.8 Humans and Robots . . . . .	968
26.9 Alternative Robotic Frameworks . . . . .	975
26.10 Application Domains . . . . .	978
Summary . . . . .	981
Bibliographical and Historical Notes . . . . .	982
<b>27 Computer Vision</b>	<b>988</b>
27.1 Introduction . . . . .	988
27.2 Image Formation . . . . .	989
27.3 Simple Image Features . . . . .	995
27.4 Classifying Images . . . . .	1002
27.5 Detecting Objects . . . . .	1006

27.6	The 3D World . . . . .	1008
27.7	Using Computer Vision . . . . .	1013
	Summary . . . . .	1026
	Bibliographical and Historical Notes . . . . .	1027

## VII Conclusions

<b>28</b>	<b>Philosophy, Ethics, and Safety of AI</b>	<b>1032</b>
28.1	The Limits of AI . . . . .	1032
28.2	Can Machines Really Think? . . . . .	1035
28.3	The Ethics of AI . . . . .	1037
	Summary . . . . .	1056
	Bibliographical and Historical Notes . . . . .	1057
<b>29</b>	<b>The Future of AI</b>	<b>1063</b>
29.1	AI Components . . . . .	1063
29.2	AI Architectures . . . . .	1069
<b>A</b>	<b>Mathematical Background</b>	<b>1074</b>
A.1	Complexity Analysis and $O()$ Notation . . . . .	1074
A.2	Vectors, Matrices, and Linear Algebra . . . . .	1076
A.3	Probability Distributions . . . . .	1078
	Bibliographical and Historical Notes . . . . .	1080
<b>B</b>	<b>Notes on Languages and Algorithms</b>	<b>1081</b>
B.1	Defining Languages with Backus–Naur Form (BNF) . . . . .	1081
B.2	Describing Algorithms with Pseudocode . . . . .	1082
B.3	Online Supplemental Material . . . . .	1083
	<b>Bibliography</b>	<b>1084</b>
	<b>Index</b>	<b>1119</b>

*This page is intentionally left blank*

# INTRODUCTION

*In which we try to explain why we consider artificial intelligence to be a subject most worthy of study, and in which we try to decide what exactly it is, this being a good thing to decide before embarking.*

We call ourselves *Homo sapiens*—man the wise—because our **intelligence** is so important to us. For thousands of years, we have tried to understand *how we think and act*—that is, how our brain, a mere handful of matter, can perceive, understand, predict, and manipulate a world far larger and more complicated than itself. The field of **artificial intelligence**, or AI, is concerned with not just understanding but also *building* intelligent entities—machines that can compute how to act effectively and safely in a wide variety of novel situations.

Intelligence

Artificial intelligence

Surveys regularly rank AI as one of the most interesting and fastest-growing fields, and it is already generating over a trillion dollars a year in revenue. AI expert Kai-Fu Lee predicts that its impact will be “more than anything in the history of mankind.” Moreover, the intellectual frontiers of AI are wide open. Whereas a student of an older science such as physics might feel that the best ideas have already been discovered by Galileo, Newton, Curie, Einstein, and the rest, AI still has many openings for full-time masterminds.

AI currently encompasses a huge variety of subfields, ranging from the general (learning, reasoning, perception, and so on) to the specific, such as playing chess, proving mathematical theorems, writing poetry, driving a car, or diagnosing diseases. AI is relevant to any intellectual task; it is truly a universal field.

## 1.1 What Is AI?

We have claimed that AI is interesting, but we have not said what it *is*. Historically, researchers have pursued several different versions of AI. Some have defined intelligence in terms of fidelity to *human* performance, while others prefer an abstract, formal definition of intelligence called **rationality**—loosely speaking, doing the “right thing.” The subject matter itself also varies: some consider intelligence to be a property of internal *thought processes* and *reasoning*, while others focus on intelligent *behavior*, an external characterization.<sup>1</sup>

Rationality

From these two dimensions—human vs. rational<sup>2</sup> and thought vs. behavior—there are four possible combinations, and there have been adherents and research programs for all

<sup>1</sup> In the public eye, there is sometimes confusion between the terms “artificial intelligence” and “machine learning.” Machine learning is a subfield of AI that studies the ability to improve performance based on experience. Some AI systems use machine learning methods to achieve competence, but some do not.

<sup>2</sup> We are not suggesting that humans are “irrational” in the dictionary sense of “deprived of normal mental clarity.” We are merely conceding that human decisions are not always mathematically perfect.

four. The methods used are necessarily different: the pursuit of human-like intelligence must be in part an empirical science related to psychology, involving observations and hypotheses about actual human behavior and thought processes; a rationalist approach, on the other hand, involves a combination of mathematics and engineering, and connects to statistics, control theory, and economics. The various groups have both disparaged and helped each other. Let us look at the four approaches in more detail.

### 1.1.1 Acting humanly: The Turing test approach

Turing test

The **Turing test**, proposed by Alan Turing (1950), was designed as a thought experiment that would sidestep the philosophical vagueness of the question “Can a machine think?” A computer passes the test if a human interrogator, after posing some written questions, cannot tell whether the written responses come from a person or from a computer. Chapter 28 discusses the details of the test and whether a computer would really be intelligent if it passed. For now, we note that programming a computer to pass a rigorously applied test provides plenty to work on. The computer would need the following capabilities:

Natural language processing  
Knowledge representation  
Automated reasoning  
Machine learning

- **natural language processing** to communicate successfully in a human language;
- **knowledge representation** to store what it knows or hears;
- **automated reasoning** to answer questions and to draw new conclusions;
- **machine learning** to adapt to new circumstances and to detect and extrapolate patterns.

Turing viewed the *physical* simulation of a person as unnecessary to demonstrate intelligence. However, other researchers have proposed a **total Turing test**, which requires interaction with objects and people in the real world. To pass the total Turing test, a robot will need

Total Turing test

- **computer vision** and speech recognition to perceive the world;
- **robotics** to manipulate objects and move about.

Computer vision

Robotics

These six disciplines compose most of AI. Yet AI researchers have devoted little effort to passing the Turing test, believing that it is more important to study the underlying principles of intelligence. The quest for “artificial flight” succeeded when engineers and inventors stopped imitating birds and started using wind tunnels and learning about aerodynamics. Aeronautical engineering texts do not define the goal of their field as making “machines that fly so exactly like pigeons that they can fool even other pigeons.”

### 1.1.2 Thinking humanly: The cognitive modeling approach

To say that a program thinks like a human, we must know how humans think. We can learn about human thought in three ways:

Introspection  
Psychological experiment  
Brain imaging

- **introspection**—trying to catch our own thoughts as they go by;
- **psychological experiments**—observing a person in action;
- **brain imaging**—observing the brain in action.

Once we have a sufficiently precise theory of the mind, it becomes possible to express the theory as a computer program. If the program’s input–output behavior matches corresponding human behavior, that is evidence that some of the program’s mechanisms could also be operating in humans.

For example, Allen Newell and Herbert Simon, who developed GPS, the “General Problem Solver” (Newell and Simon, 1961), were not content merely to have their program solve

problems correctly. They were more concerned with comparing the sequence and timing of its reasoning steps to those of human subjects solving the same problems. The interdisciplinary field of **cognitive science** brings together computer models from AI and experimental techniques from psychology to construct precise and testable theories of the human mind.

Cognitive science

Cognitive science is a fascinating field in itself, worthy of several textbooks and at least one encyclopedia (Wilson and Keil, 1999). We will occasionally comment on similarities or differences between AI techniques and human cognition. Real cognitive science, however, is necessarily based on experimental investigation of actual humans or animals. We will leave that for other books, as we assume the reader has only a computer for experimentation.

In the early days of AI there was often confusion between the approaches. An author would argue that an algorithm performs well on a task and that it is *therefore* a good model of human performance, or vice versa. Modern authors separate the two kinds of claims; this distinction has allowed both AI and cognitive science to develop more rapidly. The two fields fertilize each other, most notably in computer vision, which incorporates neurophysiological evidence into computational models. Recently, the combination of neuroimaging methods combined with machine learning techniques for analyzing such data has led to the beginnings of a capability to “read minds”—that is, to ascertain the semantic content of a person’s inner thoughts. This capability could, in turn, shed further light on how human cognition works.

### 1.1.3 Thinking rationally: The “laws of thought” approach

The Greek philosopher Aristotle was one of the first to attempt to codify “right thinking”—that is, irrefutable reasoning processes. His **sylogisms** provided patterns for argument structures that always yielded correct conclusions when given correct premises. The canonical example starts with *Socrates is a man* and *all men are mortal* and concludes that *Socrates is mortal*. (This example is probably due to Sextus Empiricus rather than Aristotle.) These laws of thought were supposed to govern the operation of the mind; their study initiated the field called **logic**.

Syllogism

Logicians in the 19th century developed a precise notation for statements about objects in the world and the relations among them. (Contrast this with ordinary arithmetic notation, which provides only for statements about *numbers*.) By 1965, programs could, in principle, solve *any* solvable problem described in logical notation. The so-called **logicist** tradition within artificial intelligence hopes to build on such programs to create intelligent systems.

Logicist

Logic as conventionally understood requires knowledge of the world that is *certain*—a condition that, in reality, is seldom achieved. We simply don’t know the rules of, say, politics or warfare in the same way that we know the rules of chess or arithmetic. The theory of **probability** fills this gap, allowing rigorous reasoning with uncertain information. In principle, it allows the construction of a comprehensive model of rational thought, leading from raw perceptual information to an understanding of how the world works to predictions about the future. What it does not do, is generate intelligent *behavior*. For that, we need a theory of rational action. Rational thought, by itself, is not enough.

Probability

### 1.1.4 Acting rationally: The rational agent approach

An **agent** is just something that acts (*agent* comes from the Latin *agere*, to do). Of course, all computer programs do something, but computer agents are expected to do more: operate autonomously, perceive their environment, persist over a prolonged time period, adapt to

Agent

Rational agent

change, and create and pursue goals. A **rational agent** is one that acts so as to achieve the best outcome or, when there is uncertainty, the best expected outcome.

In the “laws of thought” approach to AI, the emphasis was on correct inferences. Making correct inferences is sometimes *part* of being a rational agent, because one way to act rationally is to deduce that a given action is best and then to act on that conclusion. On the other hand, there are ways of acting rationally that cannot be said to involve inference. For example, recoiling from a hot stove is a reflex action that is usually more successful than a slower action taken after careful deliberation.

All the skills needed for the Turing test also allow an agent to act rationally. Knowledge representation and reasoning enable agents to reach good decisions. We need to be able to generate comprehensible sentences in natural language to get by in a complex society. We need learning not only for erudition, but also because it improves our ability to generate effective behavior, especially in circumstances that are new.

The rational-agent approach to AI has two advantages over the other approaches. First, it is more general than the “laws of thought” approach because correct inference is just one of several possible mechanisms for achieving rationality. Second, it is more amenable to scientific development. The standard of rationality is mathematically well defined and completely general. We can often work back from this specification to derive agent designs that provably achieve it—something that is largely impossible if the goal is to imitate human behavior or thought processes.

Do the right thing



For these reasons, the rational-agent approach to AI has prevailed throughout most of the field’s history. In the early decades, rational agents were built on logical foundations and formed definite plans to achieve specific goals. Later, methods based on probability theory and machine learning allowed the creation of agents that could make decisions under uncertainty to attain the best expected outcome. In a nutshell, *AI has focused on the study and construction of agents that **do the right thing***. What counts as the right thing is defined by the objective that we provide to the agent. This general paradigm is so pervasive that we might call it the **standard model**. It prevails not only in AI, but also in control theory, where a controller minimizes a cost function; in operations research, where a policy maximizes a sum of rewards; in statistics, where a decision rule minimizes a loss function; and in economics, where a decision maker maximizes utility or some measure of social welfare.

Standard model

Limited rationality

We need to make one important refinement to the standard model to account for the fact that perfect rationality—always taking the exactly optimal action—is not feasible in complex environments. The computational demands are just too high. Chapters 6 and 16 deal with the issue of **limited rationality**—acting appropriately when there is not enough time to do all the computations one might like. However, perfect rationality often remains a good starting point for theoretical analysis.

### 1.1.5 Beneficial machines

The standard model has been a useful guide for AI research since its inception, but it is probably not the right model in the long run. The reason is that the standard model assumes that we will supply a fully specified objective to the machine.

For an artificially defined task such as chess or shortest-path computation, the task comes with an objective built in—so the standard model is applicable. As we move into the real world, however, it becomes more and more difficult to specify the objective completely and

correctly. For example, in designing a self-driving car, one might think that the objective is to reach the destination safely. But driving along any road incurs a risk of injury due to other errant drivers, equipment failure, and so on; thus, a strict goal of safety requires staying in the garage. There is a tradeoff between making progress towards the destination and incurring a risk of injury. How should this tradeoff be made? Furthermore, to what extent can we allow the car to take actions that would annoy other drivers? How much should the car moderate its acceleration, steering, and braking to avoid shaking up the passenger? These kinds of questions are difficult to answer a priori. They are particularly problematic in the general area of human–robot interaction, of which the self-driving car is one example.

The problem of achieving agreement between our true preferences and the objective we put into the machine is called the **value alignment problem**: the values or objectives put into the machine must be aligned with those of the human. If we are developing an AI system in the lab or in a simulator—as has been the case for most of the field’s history—there is an easy fix for an incorrectly specified objective: reset the system, fix the objective, and try again. As the field progresses towards increasingly capable intelligent systems that are deployed in the real world, this approach is no longer viable. A system deployed with an incorrect objective will have negative consequences. Moreover, the more intelligent the system, the more negative the consequences.

Value alignment  
problem

Returning to the apparently unproblematic example of chess, consider what happens if the machine is intelligent enough to reason and act beyond the confines of the chessboard. In that case, it might attempt to increase its chances of winning by such ruses as hypnotizing or blackmailing its opponent or bribing the audience to make rustling noises during its opponent’s thinking time.<sup>3</sup> It might also attempt to hijack additional computing power for itself. *These behaviors are not “unintelligent” or “insane”; they are a logical consequence of defining winning as the sole objective for the machine.*



It is impossible to anticipate all the ways in which a machine pursuing a fixed objective might misbehave. There is good reason, then, to think that the standard model is inadequate. We don’t want machines that are intelligent in the sense of pursuing *their* objectives; we want them to pursue *our* objectives. If we cannot transfer those objectives perfectly to the machine, then we need a new formulation—one in which the machine is pursuing our objectives, but is necessarily *uncertain* as to what they are. When a machine knows that it doesn’t know the complete objective, it has an incentive to act cautiously, to ask permission, to learn more about our preferences through observation, and to defer to human control. Ultimately, we want agents that are **provably beneficial** to humans. We will return to this topic in Section 1.5.

Provably beneficial

## 1.2 The Foundations of Artificial Intelligence

In this section, we provide a brief history of the disciplines that contributed ideas, viewpoints, and techniques to AI. Like any history, this one concentrates on a small number of people, events, and ideas and ignores others that also were important. We organize the history around a series of questions. We certainly would not wish to give the impression that these questions are the only ones the disciplines address or that the disciplines have all been working toward AI as their ultimate fruition.

<sup>3</sup> In one of the first books on chess, Ruy Lopez (1561) wrote, “Always place the board so the sun is in your opponent’s eyes.”

### 1.2.1 Philosophy

- Can formal rules be used to draw valid conclusions?
- How does the mind arise from a physical brain?
- Where does knowledge come from?
- How does knowledge lead to action?

Aristotle (384–322 BCE) was the first to formulate a precise set of laws governing the rational part of the mind. He developed an informal system of syllogisms for proper reasoning, which in principle allowed one to generate conclusions mechanically, given initial premises.

Ramon Llull (c. 1232–1315) devised a system of reasoning published as *Ars Magna* or *The Great Art* (1305). Llull tried to implement his system using an actual mechanical device: a set of paper wheels that could be rotated into different permutations.

Around 1500, Leonardo da Vinci (1452–1519) designed but did not build a mechanical calculator; recent reconstructions have shown the design to be functional. The first known calculating machine was constructed around 1623 by the German scientist Wilhelm Schickard (1592–1635). Blaise Pascal (1623–1662) built the Pascaline in 1642 and wrote that it “produces effects which appear nearer to thought than all the actions of animals.” Gottfried Wilhelm Leibniz (1646–1716) built a mechanical device intended to carry out operations on concepts rather than numbers, but its scope was rather limited. In his 1651 book *Leviathan*, Thomas Hobbes (1588–1679) suggested the idea of a thinking machine, an “artificial animal” in his words, arguing “For what is the heart but a spring; and the nerves, but so many strings; and the joints, but so many wheels.” He also suggested that reasoning was like numerical computation: “For ‘reason’ . . . is nothing but ‘reckoning,’ that is adding and subtracting.”

It’s one thing to say that the mind operates, at least in part, according to logical or numerical rules, and to build physical systems that emulate some of those rules. It’s another to say that the mind itself *is* such a physical system. René Descartes (1596–1650) gave the first clear discussion of the distinction between mind and matter. He noted that a purely physical conception of the mind seems to leave little room for free will. If the mind is governed entirely by physical laws, then it has no more free will than a rock “deciding” to fall downward. Descartes was a proponent of **dualism**. He held that there is a part of the human mind (or soul or spirit) that is outside of nature, exempt from physical laws. Animals, on the other hand, did not possess this dual quality; they could be treated as machines.

An alternative to dualism is **materialism**, which holds that the brain’s operation according to the laws of physics *constitutes* the mind. Free will is simply the way that the perception of available choices appears to the choosing entity. The terms **physicalism** and **naturalism** are also used to describe this view that stands in contrast to the supernatural.

Given a physical mind that manipulates knowledge, the next problem is to establish the source of knowledge. The **empiricism** movement, starting with Francis Bacon’s (1561–1626) *Novum Organum*,<sup>4</sup> is characterized by a dictum of John Locke (1632–1704): “Nothing is in the understanding, which was not first in the senses.”

David Hume’s (1711–1776) *A Treatise of Human Nature* (Hume, 1739) proposed what is now known as the principle of **induction**: that general rules are acquired by exposure to repeated associations between their elements.

<sup>4</sup> The *Novum Organum* is an update of Aristotle’s *Organon*, or instrument of thought.

Dualism

Empiricism

Induction

Building on the work of Ludwig Wittgenstein (1889–1951) and Bertrand Russell (1872–1970), the famous Vienna Circle (Sigmund, 2017), a group of philosophers and mathematicians meeting in Vienna in the 1920s and 1930s, developed the doctrine of **logical positivism**. This doctrine holds that all knowledge can be characterized by logical theories connected, ultimately, to **observation sentences** that correspond to sensory inputs; thus logical positivism combines rationalism and empiricism.

Logical positivism

Observation sentence

The **confirmation theory** of Rudolf Carnap (1891–1970) and Carl Hempel (1905–1997) attempted to analyze the acquisition of knowledge from experience by quantifying the degree of belief that should be assigned to logical sentences based on their connection to observations that confirm or disconfirm them. Carnap’s book *The Logical Structure of the World* (1928) was perhaps the first theory of mind as a computational process.

Confirmation theory

The final element in the philosophical picture of the mind is the connection between knowledge and action. This question is vital to AI because intelligence requires action as well as reasoning. Moreover, only by understanding how actions are justified can we understand how to build an agent whose actions are justifiable (or rational).

Aristotle argued (in *De Motu Animalium*) that actions are justified by a logical connection between goals and knowledge of the action’s outcome:

But how does it happen that thinking is sometimes accompanied by action and sometimes not, sometimes by motion, and sometimes not? It looks as if almost the same thing happens as in the case of reasoning and making inferences about unchanging objects. But in that case the end is a speculative proposition . . . whereas here the conclusion which results from the two premises is an action. . . . I need covering; a cloak is a covering. I need a cloak. What I need, I have to make; I need a cloak. I have to make a cloak. And the conclusion, the “I have to make a cloak,” is an action.

In the *Nicomachean Ethics* (Book III. 3, 1112b), Aristotle further elaborates on this topic, suggesting an algorithm:

We deliberate not about ends, but about means. For a doctor does not deliberate whether he shall heal, nor an orator whether he shall persuade, . . . They assume the end and consider how and by what means it is attained, and if it seems easily and best produced thereby; while if it is achieved by one means only they consider *how* it will be achieved by this and by what means *this* will be achieved, till they come to the first cause, . . . and what is last in the order of analysis seems to be first in the order of becoming. And if we come on an impossibility, we give up the search, e.g., if we need money and this cannot be got; but if a thing appears possible we try to do it.

Aristotle’s algorithm was implemented 2300 years later by Newell and Simon in their **General Problem Solver** program. We would now call it a greedy regression planning system (see Chapter 11). Methods based on logical planning to achieve definite goals dominated the first few decades of theoretical research in AI.

Thinking purely in terms of actions achieving goals is often useful but sometimes inapplicable. For example, if there are several different ways to achieve a goal, there needs to be some way to choose among them. More importantly, it may not be possible to achieve a goal with certainty, but some action must still be taken. How then should one decide? Antoine Arnauld (1662), analyzing the notion of rational decisions in gambling, proposed a quantitative formula for maximizing the expected monetary value of the outcome. Later, Daniel Bernoulli (1738) introduced the more general notion of **utility** to capture the internal, subjective value

Utility

of an outcome. The modern notion of rational decision making under uncertainty involves maximizing expected utility, as explained in Chapter 15.

#### Utilitarianism

In matters of ethics and public policy, a decision maker must consider the interests of multiple individuals. Jeremy Bentham (1823) and John Stuart Mill (1863) promoted the idea of **utilitarianism**: that rational decision making based on maximizing utility should apply to all spheres of human activity, including public policy decisions made on behalf of many individuals. Utilitarianism is a specific kind of **consequentialism**: the idea that what is right and wrong is determined by the expected outcomes of an action.

#### Deontological ethics

In contrast, Immanuel Kant, in 1785, proposed a theory of rule-based or **deontological ethics**, in which “doing the right thing” is determined not by outcomes but by universal social laws that govern allowable actions, such as “don’t lie” or “don’t kill.” Thus, a utilitarian could tell a white lie if the expected good outweighs the bad, but a Kantian would be bound not to, because lying is inherently wrong. Mill acknowledged the value of rules, but understood them as efficient decision procedures compiled from first-principles reasoning about consequences. Many modern AI systems adopt exactly this approach.

### 1.2.2 Mathematics

- What are the formal rules to draw valid conclusions?
- What can be computed?
- How do we reason with uncertain information?

Philosophers staked out some of the fundamental ideas of AI, but the leap to a formal science required the mathematization of logic and probability and the introduction of a new branch of mathematics: computation.

#### Formal logic

The idea of **formal logic** can be traced back to the philosophers of ancient Greece, India, and China, but its mathematical development really began with the work of George Boole (1815–1864), who worked out the details of propositional, or Boolean, logic (Boole, 1847). In 1879, Gottlob Frege (1848–1925) extended Boole’s logic to include objects and relations, creating the first-order logic that is used today.<sup>5</sup> In addition to its central role in the early period of AI research, first-order logic motivated the work of Gödel and Turing that underpinned computation itself, as we explain below.

#### Probability

The theory of **probability** can be seen as generalizing logic to situations with uncertain information—a consideration of great importance for AI. Gerolamo Cardano (1501–1576) first framed the idea of probability, describing it in terms of the possible outcomes of gambling events. In 1654, Blaise Pascal (1623–1662), in a letter to Pierre Fermat (1601–1665), showed how to predict the future of an unfinished gambling game and assign average pay-offs to the gamblers. Probability quickly became an invaluable part of the quantitative sciences, helping to deal with uncertain measurements and incomplete theories. Jacob Bernoulli (1654–1705, uncle of Daniel), Pierre Laplace (1749–1827), and others advanced the theory and introduced new statistical methods. Thomas Bayes (1702–1761) proposed a rule for updating probabilities in the light of new evidence; Bayes’ rule is a crucial tool for AI systems.

#### Statistics

The formalization of probability, combined with the availability of data, led to the emergence of **statistics** as a field. One of the first uses was John Graunt’s analysis of Lon-

<sup>5</sup> Frege’s proposed notation for first-order logic—an arcane combination of textual and geometric features—never became popular.

don census data in 1662. Ronald Fisher is considered the first modern statistician (Fisher, 1922). He brought together the ideas of probability, experiment design, analysis of data, and computing—in 1919, he insisted that he couldn’t do his work without a mechanical calculator called the MILLIONAIRE (the first calculator that could do multiplication), even though the cost of the calculator was more than his annual salary (Ross, 2012).

The history of computation is as old as the history of numbers, but the first nontrivial **algorithm** is thought to be Euclid’s algorithm for computing greatest common divisors. The word *algorithm* comes from Muhammad ibn Musa al-Khwarizmi, a 9th century mathematician, whose writings also introduced Arabic numerals and algebra to Europe. Boole and others discussed algorithms for logical deduction, and, by the late 19th century, efforts were under way to formalize general mathematical reasoning as logical deduction.

Algorithm

Kurt Gödel (1906–1978) showed that there exists an effective procedure to prove any true statement in the first-order logic of Frege and Russell, but that first-order logic could not capture the principle of mathematical induction needed to characterize the natural numbers. In 1931, Gödel showed that limits on deduction do exist. His **incompleteness theorem** showed that in any formal theory as strong as Peano arithmetic (the elementary theory of natural numbers), there are necessarily true statements that have no proof within the theory.

Incompleteness theorem

This fundamental result can also be interpreted as showing that some functions on the integers cannot be represented by an algorithm—that is, they cannot be computed. This motivated Alan Turing (1912–1954) to try to characterize exactly which functions *are* **computable**—capable of being computed by an effective procedure. The Church–Turing thesis proposes to identify the general notion of computability with functions computed by a Turing machine (Turing, 1936). Turing also showed that there were some functions that no Turing machine can compute. For example, no machine can tell *in general* whether a given program will return an answer on a given input or run forever.

Computability

Although computability is important to an understanding of computation, the notion of **tractability** has had an even greater impact on AI. Roughly speaking, a problem is called intractable if the time required to solve instances of the problem grows exponentially with the size of the instances. The distinction between polynomial and exponential growth in complexity was first emphasized in the mid-1960s (Cobham, 1964; Edmonds, 1965). It is important because exponential growth means that even moderately large instances cannot be solved in any reasonable time.

Tractability

The theory of **NP-completeness**, pioneered by Cook (1971) and Karp (1972), provides a basis for analyzing the tractability of problems: any problem class to which the class of NP-complete problems can be reduced is likely to be intractable. (Although it has not been proved that NP-complete problems are necessarily intractable, most theoreticians believe it.) These results contrast with the optimism with which the popular press greeted the first computers—“Electronic Super-Brains” that were “Faster than Einstein!” Despite the increasing speed of computers, careful use of resources and necessary imperfection will characterize intelligent systems. Put crudely, the world is an *extremely* large problem instance!

NP-completeness

### 1.2.3 Economics

- How should we make decisions in accordance with our preferences?
- How should we do this when others may not go along?
- How should we do this when the payoff may be far in the future?

The science of economics originated in 1776, when Adam Smith (1723–1790) published *An Inquiry into the Nature and Causes of the Wealth of Nations*. Smith proposed to analyze economies as consisting of many individual agents attending to their own interests. Smith was not, however, advocating financial greed as a moral position: his earlier (1759) book *The Theory of Moral Sentiments* begins by pointing out that concern for the well-being of others is an essential component of the interests of every individual.

Most people think of economics as being about money, and indeed the first mathematical analysis of decisions under uncertainty, the maximum-expected-value formula of Arnauld (1662), dealt with the monetary value of bets. Daniel Bernoulli (1738) noticed that this formula didn't seem to work well for larger amounts of money, such as investments in maritime trading expeditions. He proposed instead a principle based on maximization of expected utility, and explained human investment choices by proposing that the marginal utility of an additional quantity of money diminished as one acquired more money.

Léon Walras (pronounced “Valrasse”) (1834–1910) gave utility theory a more general foundation in terms of preferences between gambles on any outcomes (not just monetary outcomes). The theory was improved by Ramsey (1931) and later by John von Neumann and Oskar Morgenstern in their book *The Theory of Games and Economic Behavior* (1944). Economics is no longer the study of money; rather it is the study of desires and preferences.

Decision theory

**Decision theory**, which combines probability theory with utility theory, provides a formal and complete framework for individual decisions (economic or otherwise) made under uncertainty—that is, in cases where probabilistic descriptions appropriately capture the decision maker's environment. This is suitable for “large” economies where each agent need pay no attention to the actions of other agents as individuals. For “small” economies, the situation is much more like a **game**: the actions of one player can significantly affect the utility of another (either positively or negatively). Von Neumann and Morgenstern's development of **game theory** (see also Luce and Raiffa, 1957) included the surprising result that, for some games, a rational agent should adopt policies that are (or least appear to be) randomized. Unlike decision theory, game theory does not offer an unambiguous prescription for selecting actions. In AI, decisions involving multiple agents are studied under the heading of **multiagent systems** (Chapter 17).

Operations research

Economists, with some exceptions, did not address the third question listed above: how to make rational decisions when payoffs from actions are not immediate but instead result from several actions taken *in sequence*. This topic was pursued in the field of **operations research**, which emerged in World War II from efforts in Britain to optimize radar installations, and later found innumerable civilian applications. The work of Richard Bellman (1957) formalized a class of sequential decision problems called **Markov decision processes**, which we study in Chapter 16 and, under the heading of **reinforcement learning**, in Chapter 23.

Satisficing

Work in economics and operations research has contributed much to our notion of rational agents, yet for many years AI research developed along entirely separate paths. One reason was the apparent complexity of making rational decisions. The pioneering AI researcher Herbert Simon (1916–2001) won the Nobel Prize in economics in 1978 for his early work showing that models based on **satisficing**—making decisions that are “good enough,” rather than laboriously calculating an optimal decision—gave a better description of actual human behavior (Simon, 1947). Since the 1990s, there has been a resurgence of interest in decision-theoretic techniques for AI.

### 1.2.4 Neuroscience

- How do brains process information?

**Neuroscience** is the study of the nervous system, particularly the brain. Although the exact way in which the brain enables thought is one of the great mysteries of science, the fact that it *does* enable thought has been appreciated for thousands of years because of the evidence that strong blows to the head can lead to mental incapacitation. It has also long been known that human brains are somehow different; in about 335 BCE Aristotle wrote, “Of all the animals, man has the largest brain in proportion to his size.”<sup>6</sup> Still, it was not until the middle of the 18th century that the brain was widely recognized as the seat of consciousness. Before then, candidate locations included the heart and the spleen.

Paul Broca’s (1824–1880) investigation of aphasia (speech deficit) in brain-damaged patients in 1861 initiated the study of the brain’s functional organization by identifying a localized area in the left hemisphere—now called Broca’s area—that is responsible for speech production.<sup>7</sup> By that time, it was known that the brain consisted largely of nerve cells, or **neurons**, but it was not until 1873 that Camillo Golgi (1843–1926) developed a staining technique allowing the observation of individual neurons (see Figure 1.1). This technique was used by Santiago Ramon y Cajal (1852–1934) in his pioneering studies of neuronal organization.<sup>8</sup> It is now widely accepted that cognitive functions result from the electrochemical operation of these structures. That is, *a collection of simple cells can lead to thought, action, and consciousness*. In the pithy words of John Searle (1992), *brains cause minds*.

We now have some data on the mapping between areas of the brain and the parts of the body that they control or from which they receive sensory input. Such mappings are able to change radically over the course of a few weeks, and some animals seem to have multiple maps. Moreover, we do not fully understand how other areas can take over functions when one area is damaged. There is almost no theory on how an individual memory is stored or on how higher-level cognitive functions operate.

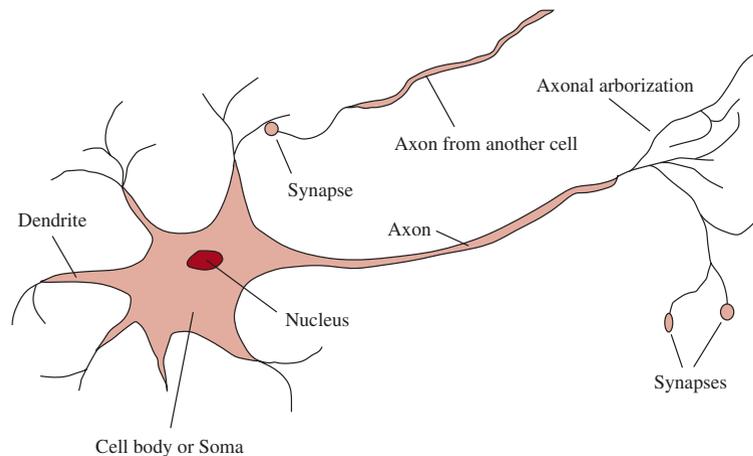
The measurement of intact brain activity began in 1929 with the invention by Hans Berger of the electroencephalograph (EEG). The development of functional magnetic resonance imaging (fMRI) (Ogawa *et al.*, 1990; Cabeza and Nyberg, 2001) is giving neuroscientists unprecedentedly detailed images of brain activity, enabling measurements that correspond in interesting ways to ongoing cognitive processes. These are augmented by advances in single-cell electrical recording of neuron activity and by the methods of **optogenetics** (Crick, 1999; Zemelman *et al.*, 2002; Han and Boyden, 2007), which allow both measurement and control of individual neurons modified to be light-sensitive.

The development of **brain–machine interfaces** (Lebedev and Nicolelis, 2006) for both sensing and motor control not only promises to restore function to disabled individuals, but also sheds light on many aspects of neural systems. A remarkable finding from this work is that the brain is able to adjust itself to interface successfully with an external device, treating it in effect like another sensory organ or limb.

<sup>6</sup> It has since been discovered that the tree shrew and some bird species exceed the human brain/body ratio.

<sup>7</sup> Many cite Alexander Hood (1824) as a possible prior source.

<sup>8</sup> Golgi persisted in his belief that the brain’s functions were carried out primarily in a continuous medium in which neurons were embedded, whereas Cajal propounded the “neuronal doctrine.” The two shared the Nobel Prize in 1906 but gave mutually antagonistic acceptance speeches.



**Figure 1.1** The parts of a nerve cell or neuron. Each neuron consists of a cell body, or soma, that contains a cell nucleus. Branching out from the cell body are a number of fibers called dendrites and a single long fiber called the axon. The axon stretches out for a long distance, much longer than the scale in this diagram indicates. Typically, an axon is 1 cm long (100 times the diameter of the cell body), but can reach up to 1 meter. A neuron makes connections with 10 to 100,000 other neurons at junctions called synapses. Signals are propagated from neuron to neuron by a complicated electrochemical reaction. The signals control brain activity in the short term and also enable long-term changes in the connectivity of neurons. These mechanisms are thought to form the basis for learning in the brain. Most information processing goes on in the cerebral cortex, the outer layer of the brain. The basic organizational unit appears to be a column of tissue about 0.5 mm in diameter, containing about 20,000 neurons and extending the full depth of the cortex (about 4 mm in humans).

Brains and digital computers have somewhat different properties. Figure 1.2 shows that computers have a cycle time that is a million times faster than a brain. The brain makes up for that with far more storage and interconnection than even a high-end personal computer, although the largest supercomputers match the brain on some metrics. Futurists make much of these numbers, pointing to an approaching **singularity** at which computers reach a superhuman level of performance (Vinge, 1993; Kurzweil, 2005; Doctorow and Stross, 2012), and then rapidly improve themselves even further. But the comparisons of raw numbers are not especially informative. Even with a computer of virtually unlimited capacity, we still require further conceptual breakthroughs in our understanding of intelligence (see Chapter 29). Crudely put, without the right theory, faster machines just give you the wrong answer faster.

Singularity

### 1.2.5 Psychology

- How do humans and animals think and act?

The origins of scientific psychology are usually traced to the work of the German physicist Hermann von Helmholtz (1821–1894) and his student Wilhelm Wundt (1832–1920). Helmholtz applied the scientific method to the study of human vision, and his *Handbook of Physiological Optics* has been described as “the single most important treatise on the physics and physiology of human vision” (Nalwa, 1993, p.15). In 1879, Wundt opened the first laboratory of experimental psychology, at the University of Leipzig. Wundt insisted on carefully