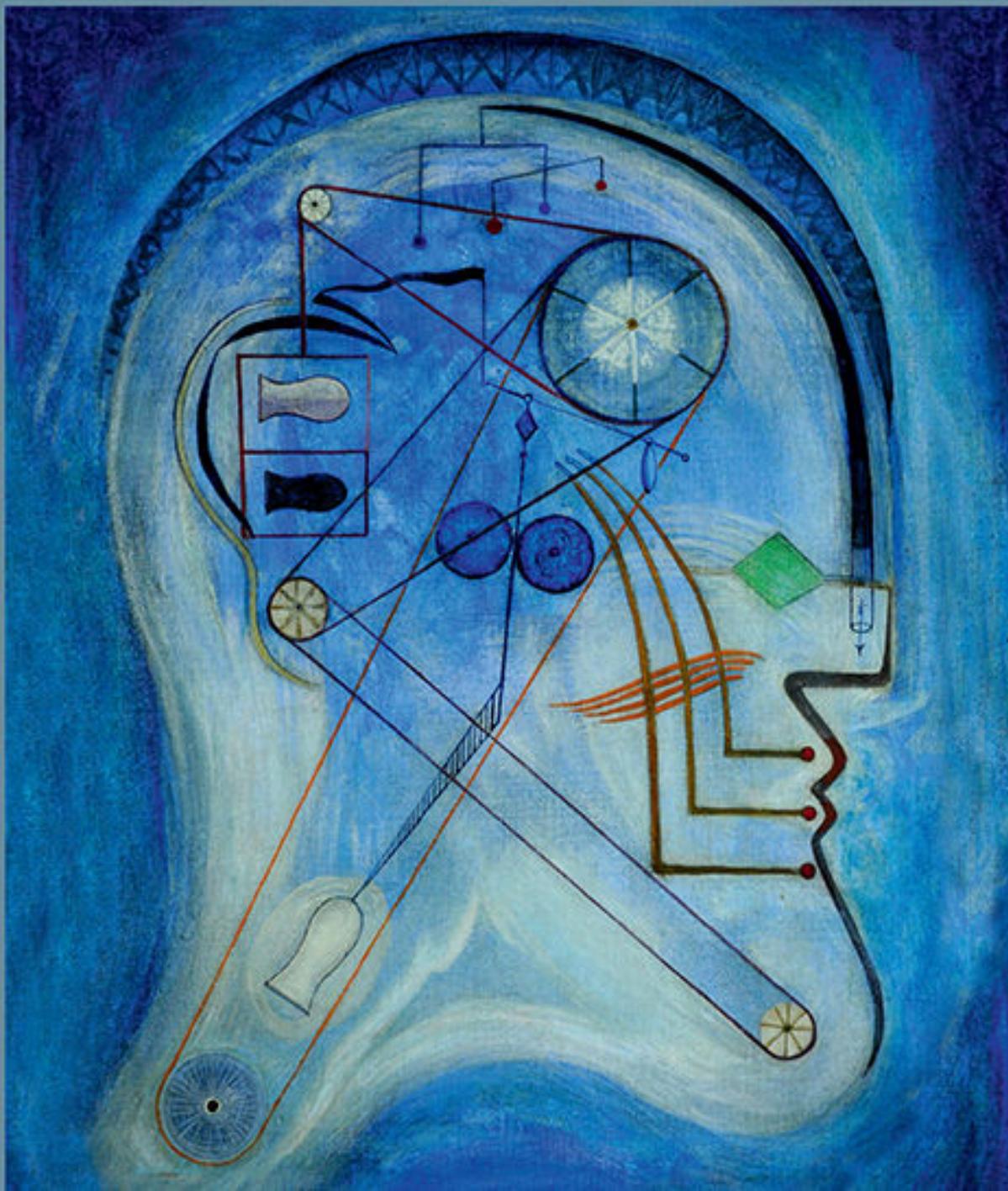


José Luis Bermúdez

# Cognitive Science

An Introduction to the Science of the Mind

Third Edition



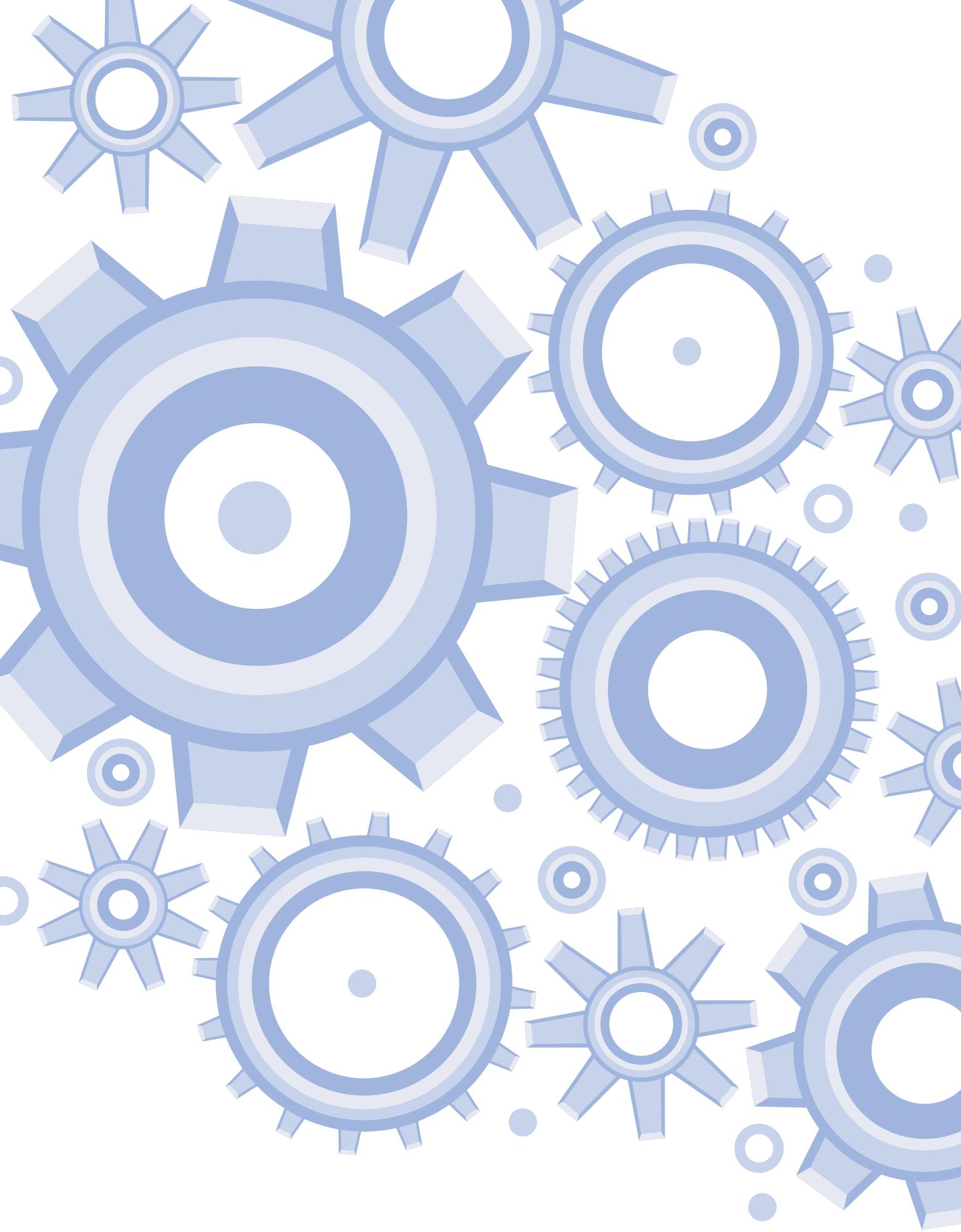




## COGNITIVE SCIENCE

### Third edition

The third edition of this popular and engaging text consolidates the interdisciplinary streams of cognitive science to present a unified narrative of cognitive science as a discipline in its own right. It teaches students to apply the techniques and theories of the cognitive scientist's "tool kit" – the vast range of methods and tools that cognitive scientists use to study the mind. Thematically organized, *Cognitive Science* underscores the problems and solutions of cognitive science rather than more narrowly examining individually the subjects that contribute to it – psychology, neuroscience, linguistics, and so on. The generous use of examples, illustrations, and applications demonstrates how theory is applied to unlock the mysteries of the human mind. Drawing upon cutting-edge research, the text has been substantially revised, with new material on Bayesian approaches to the mind and on deep learning. An extensive online set of resources is available to aid instructors and students alike. Sample syllabi show how the text can support a variety of courses, making it a highly flexible teaching and learning resource at both the undergraduate and graduate levels.



# COGNITIVE SCIENCE

An Introduction to the Science of the Mind

Third edition

José Luis Bermúdez  
Texas A&M University



University Printing House, Cambridge CB2 8BS, United Kingdom  
One Liberty Plaza, 20th Floor, New York, NY 10006, USA  
477 Williamstown Road, Port Melbourne, VIC 3207, Australia  
314–321, 3rd Floor, Plot 3, Splendor Forum, Jasola District Centre, New Delhi – 110025, India  
79 Anson Road, #06–04/06, Singapore 079906

Cambridge University Press is part of the University of Cambridge.

It furthers the University's mission by disseminating knowledge in the pursuit of education, learning, and research at the highest international levels of excellence.

[www.cambridge.org](http://www.cambridge.org)

Information on this title: [www.cambridge.org/9781108424493](http://www.cambridge.org/9781108424493)

DOI: [10.1017/9781108339216](https://doi.org/10.1017/9781108339216)

© Cambridge University Press 2020

This publication is in copyright. Subject to statutory exception and to the provisions of relevant collective licensing agreements, no reproduction of any part may take place without the written permission of Cambridge University Press.

First published 2020

Printed in Singapore by Markono Print Media Pte Ltd

*A catalogue record for this publication is available from the British Library.*

ISBN 978-1-108-42449-3 Hardback

ISBN 978-1-108-44034-9 Paperback

Additional resources for this publication at [www.cambridge.org/bermudez3e](http://www.cambridge.org/bermudez3e)

Cambridge University Press has no responsibility for the persistence or accuracy of URLs for external or third-party internet websites referred to in this publication and does not guarantee that any content on such websites is, or will remain, accurate or appropriate.



## CONTENTS

*List of Boxes* *xiv*

*List of Figures* *xv*

*List of Tables* *xxii*

*Preface* *xxiii*

*Acknowledgments for the First Edition* *xxxii*

*Acknowledgments for the Second Edition* *xxxii*

*Acknowledgments for the Third Edition* *xxxiii*

Introduction: The Challenge of Cognitive Science *3*

### **PART I HISTORICAL LANDMARKS *12***

- 1** The Prehistory of Cognitive Science *15*
- 2** The Discipline Matures: Three Milestones *37*
- 3** The Turn to the Brain *65*

### **PART II MODELS AND TOOLS *96***

- 4** Physical Symbol Systems and the Language of Thought *99*
- 5** Neural Networks and Distributed Information Processing *123*
- 6** Applying Dynamical Systems Theory to Model the Mind *149*
- 7** Bayesianism in Cognitive Science *171*
- 8** Modules and Architectures *203*
- 9** Strategies for Brain Mapping *229*

### **PART III APPLICATIONS *256***

- 10** Models of Language Learning *259*
- 11** Object Perception and Folk Physics *285*
- 12** Machine Learning: From Expert Systems to Deep Learning *307*
- 13** Exploring Mindreading *335*
- 14** Mindreading: Advanced Topics *357*
- 15** The Cognitive Science of Consciousness *379*
- 16** Robotics: From GOFAI to Situated Cognition and Behavior-Based Robotics *407*
- 17** Looking Ahead: Challenges and Opportunities *437*

*Glossary* *444*

*Bibliography* *454*

*Index for Cognitive Science (3rd edition)* *478*





## CONTENTS

*List of Boxes* xiv

*List of Figures* xv

*List of Tables* xxii

*Preface* xxiii

*Acknowledgments for the First Edition* xxxi

*Acknowledgments for the Second Edition* xxxii

*Acknowledgments for the Third Edition* xxxiii

**Introduction: The Challenge of Cognitive Science** 3

**0.1 Cognitive Science: An Interdisciplinary Endeavor** 3

**0.2 Levels of Explanation: The Contrast between Psychology and Neuroscience** 5

How Psychology Is Organized 6

How Neuroscience Is Organized 7

**0.3 The Challenge of Cognitive Science** 10

Three Dimensions of Variation 10

The Space of Cognitive Science 10

## PART I HISTORICAL LANDMARKS 12

1 **The Prehistory of Cognitive Science** 15

**1.1 The Reaction against Behaviorism in Psychology** 16

Learning without Reinforcement: Tolman and Honzik, "Insight" in Rats" (1930) 17

Cognitive Maps in Rats? Tolman, Ritchie, and Kalish, "Studies in Spatial Learning" (1946) 20

Plans and Complex Behaviors: Lashley, "The Problem of Serial Order in Behavior" (1951) 21

**1.2 The Theory of Computation and the Idea of an Algorithm** 22

Algorithms and Turing Machines: Turing, "On Computable Numbers, with an Application to the Decision Problem" (1936–7) 23

**1.3 Linguistics and the Formal Analysis of Language** 25

The Structure of Language: Chomsky's *Syntactic Structures* (1957) 26

<b>1.4 Information-Processing Models in Psychology</b>	<b>28</b>
How Much Information Can We Handle? George Miller's "The Magical Number Seven, Plus or Minus Two" (1956)	29
The Flow of Information: Donald Broadbent's "The Role of Auditory Localization in Attention and Memory Span" (1954) and <i>Perception and Communication</i> (1958)	30
<b>1.5 Connections and Points of Contact</b>	<b>32</b>
<b>2 The Discipline Matures: Three Milestones</b>	<b>37</b>
<b>2.1 Language and Micro-worlds</b>	<b>38</b>
Natural Language Processing: Winograd, <i>Understanding Natural Language</i> (1972)	39
SHRDLU in Action	41
<b>2.2 How Do Mental Images Represent?</b>	<b>47</b>
Mental Rotation: Shepard and Metzler, "Mental Rotation of Three-Dimensional Objects" (1971)	48
Information Processing in Mental Imagery	50
<b>2.3 An Interdisciplinary Model of Vision</b>	<b>53</b>
Levels of Explanation: Marr's <i>Vision</i> (1982)	53
Applying Top-Down Analysis to the Visual System	55
<b>3 The Turn to the Brain</b>	<b>65</b>
<b>3.1 Cognitive Systems as Functional Systems?</b>	<b>66</b>
<b>3.2 The Anatomy of the Brain and the Primary Visual Pathway</b>	<b>68</b>
The Two Visual Systems Hypothesis: Ungerleider and Mishkin, "Two Cortical Visual Systems" (1982)	70
<b>3.3 Extending Computational Modeling to the Brain</b>	<b>76</b>
A New Set of Algorithms: Rumelhart, McClelland, and the PDP Research Group, <i>Parallel Distributed Processing: Explorations in the Microstructure of Cognition</i> (1986)	77
Pattern Recognition in Neural Networks: Gorman and Sejnowski, "Analysis of Hidden Units in a Layered Network Trained to Identify Sonar Targets" (1998)	78
<b>3.4 Mapping the Stages of Lexical Processing</b>	<b>80</b>
Functional Neuroimaging with PET: Petersen, Fox, Posner, and Mintun, "Positron Emission Tomographic Studies of the Cortical Anatomy of Single-Word Processing" (1988)	81
Petersen, Fox, Posner, and Mintun, "Positron Emission Tomographic Studies of the Cortical Anatomy of Single-Word Processing" (1988)	81
<b>3.5 Studying Memory for Visual Events</b>	<b>84</b>
Functional Neuroimaging with fMRI	86
Brewer, Zhao, Desmond, Glover, and Gabrieli, "Making Memories: Brain Activity That Predicts How Well Visual Experience Will Be Remembered" (1998)	87

**3.6 The Neural Correlates of the BOLD Signal 90**

Logothetis, "The Underpinnings of the BOLD Functional Magnetic Resonance Imaging Signal" (2001) 91

**PART II MODELS AND TOOLS 96****4 Physical Symbol Systems and the Language of Thought 99****4.1 The Physical Symbol System Hypothesis 100**

Symbols and Symbol Systems 101

Transforming Symbol Structures 102

Intelligent Action and the Physical Symbol System 106

**4.2 From Physical Symbol Systems to the Language of Thought 106**

Intentional Realism and Causation by Content 108

The Language of Thought and the Relation between Syntax and Semantics 110

**4.3 The Russian Room Argument and the Turing Test 114**

Responding to the Russian Room Argument 117

**5 Neural Networks and Distributed Information Processing 123****5.1 Neurally Inspired Models of Information Processing 124**

Neurons and Network Units 125

**5.2 Single-Layer Networks and Boolean Functions 128**

Learning in Single-Layer Networks: The Perceptron Convergence Rule 131

Linear Separability and the Limits of Perceptron Convergence 134

**5.3 Multilayer Networks 137**

The Backpropagation Algorithm 138

How Biologically Plausible Are Neural Networks? 139

**5.4 Information Processing in Neural Networks: Key Features 141**

Distributed Representations 141

No Clear Distinction between Information Storage and Information Processing 142

The Ability to Learn from "Experience" 143

**6 Applying Dynamical Systems Theory to Model the Mind 149****6.1 Cognitive Science and Dynamical Systems 149**

What Are Dynamical Systems? 150

The Dynamical Systems Hypothesis: Cognitive Science without Representations? 153

**6.2 Applying Dynamical Systems: Two Examples from Child Development 158**

Two Ways of Thinking about Motor Control 159

Dynamical Systems and the A-Not-B Error 161

Assessing the Dynamical Systems Approach 166

<b>7 Bayesianism in Cognitive Science</b>	<b>171</b>
<b>    7.1 Bayesianism: A Primer</b>	<b>172</b>
Degrees of Belief and Subjective Probability	173
Conditional Probability	175
Bayes's Rule (the Short Version)	176
<b>    7.2 Perception as a Bayesian Problem</b>	<b>179</b>
The Predictive Challenge of Perception	179
Case Study: Binocular Rivalry	182
<b>    7.3 Neuroeconomics: Bayes in the Brain</b>	<b>186</b>
What Is Expected Utility?	187
Case Study: Neurons That Code for Expected Utility	190
<i>Probability-Detecting Neurons</i>	193
<i>Utility-Detecting Neurons</i>	194
<i>Combining Probability and Utility</i>	196
<b>8 Modules and Architectures</b>	<b>203</b>
<b>    8.1 Architectures for Artificial Agents</b>	<b>204</b>
Three Agent Architectures	204
<b>    8.2 Fodor on the Modularity of Mind</b>	<b>208</b>
Modular and Nonmodular Processing	208
<b>    8.3 The Massive Modularity Hypothesis</b>	<b>210</b>
The Cheater Detection Module	211
The Evolution of Cooperation	213
Two Arguments	216
Evaluating the Arguments for Massive Modularity	218
<b>    8.4 Hybrid Architectures: The Example of ACT-R</b>	<b>219</b>
The ACT-R Architecture	220
ACT-R as a Hybrid Architecture	222
<b>9 Strategies for Brain Mapping</b>	<b>229</b>
<b>    9.1 Structure and Function in the Brain</b>	<b>230</b>
Exploring Anatomical Connectivity	232
<b>    9.2 Studying Cognitive Functioning: Techniques from Neuroscience</b>	<b>237</b>
Mapping the Brain's Electrical Activity: EEG and MEG	237
Mapping the Brain's Blood Flow and Blood Oxygen Levels: PET and fMRI	240
<b>    9.3 Combining Resources I: The Locus of Selection Problem</b>	<b>241</b>
Combining ERPs and Single-Unit Recordings	242

**9.4 Combining Resources II: Networks for Attention 246**

Two Hypotheses about Visuospatial Attention 248

**9.5 From Data to Maps: Problems and Pitfalls 249**

From Blood Flow to Cognition? 250

Noise in the System? 251

Functional Connectivity versus Effective Connectivity 252

**PART III APPLICATIONS 256****10 Models of Language Learning 259****10.1 Language and Rules 260**

Understanding a Language and Learning a Language 261

**10.2 Language Learning and the Language of Thought: Fodor's Argument 263****10.3 Language Learning in Neural Networks 266**

The Challenge of Tense Learning 267

Connectionist Models of Tense Learning 269

**10.4 Bayesian Language Learning 274**

Probabilities in Word and Phrase Segmentation 275

Understanding Pronouns 276

Learning Linguistic Categories 278

**11 Object Perception and Folk Physics 285****11.1 Object Permanence and Physical Reasoning in Infancy 286**

Infant Cognition and the Dishabituation Paradigm 286

How Should the Dishabituation Experiments Be Interpreted? 292

**11.2 Neural Network Models of Children's Physical Reasoning 293**

Modeling Object Permanence 295

Modeling the Balance Beam Problem 297

**11.3 Conclusion: The Question of Levels 300****12 Machine Learning: From Expert Systems to Deep Learning 307****12.1 Expert Systems and Machine Learning 308**

Expert Systems and Decision Trees 308

ID3: An Algorithm for Machine Learning 310

**12.2 Representation Learning and Deep Learning 315**

Deep Learning and the Visual Cortex 318

**12.3 The Machinery of Deep Learning 321**

Autoencoders 322

	Convolutional Neural Networks	324
	<i>Sparse Connectivity</i>	325
	<i>Shared Weights</i>	326
	<i>Invariance under Translation</i>	326
	<b>12.4 Deep Reinforcement Learning</b>	327
13	<b>Exploring Mindreading</b>	335
	<b>13.1 Pretend Play and Metarepresentation</b>	336
	The Significance of Pretend Play	336
	Leslie on Pretend Play and Metarepresentation	337
	The Link to Mindreading	341
	<b>13.2 Metarepresentation, Autism, and Theory of Mind</b>	341
	Using the False Belief Task to Study Mindreading	342
	Interpreting the Results	344
	Implicit and Explicit Understanding of False Belief	347
	<b>13.3 The Mindreading System</b>	348
	First Steps in Mindreading	349
	From Dyadic to Triadic Interactions: Joint Visual Attention	351
	TESS and TOMM	352
14	<b>Mindreading: Advanced Topics</b>	357
	<b>14.1 Why Does It Take Children So Long to Learn to Understand False Belief?</b>	358
	Leslie's Answer: The Selection Processor Hypothesis	358
	An Alternative Model of Theory of Mind Development	360
	<b>14.2 Mindreading as Simulation</b>	363
	Standard Simulationism	363
	Radical Simulationism	365
	<b>14.3 The Cognitive Neuroscience of Mindreading</b>	365
	Neuroimaging Evidence for a Dedicated Theory of Mind System?	366
	Neuroscientific Evidence for Simulation in Low-Level Mindreading?	369
	Neuroscientific Evidence for Simulation in High-Level Mindreading?	373
15	<b>The Cognitive Science of Consciousness</b>	379
	<b>15.1 The Challenge of Consciousness: The Knowledge Argument</b>	380
	<b>15.2 Information Processing without Conscious Awareness: Some Basic Data</b>	382
	Consciousness and Priming	382
	Nonconscious Processing in Blindsight and Unilateral Spatial Neglect	384

---

<b>15.3 So What Is Consciousness For?</b>	387
What Is Missing in Blindsight and Spatial Neglect	389
Milner and Goodale: Vision for Action and Vision for Perception	389
What Is Missing in Masked Priming	392
<b>15.4 Two Types of Consciousness and the Hard Problem</b>	393
<b>15.5 The Global Workspace Theory of Consciousness</b>	396
The Building Blocks of Global Workspace Theory	396
The Global Neuronal Workspace Theory	397
<b>15.6 Conclusion</b>	400
<b>16 Robotics: From GOFAI to Situated Cognition and Behavior-Based Robotics</b>	407
<b>16.1 GOFAI Robotics: SHAKEY</b>	408
SHAKEY's Software I: Low-Level Activities and Intermediate-Level Actions	409
SHAKEY's Software II: Logic Programming in STRIPS and PLANEX	413
<b>16.2 Situated Cognition and Biorobotics</b>	414
The Challenge of Building a Situated Agent	415
Situated Cognition and Knowledge Representation	416
Biorobotics: Insects and Morphological Computation	418
<b>16.3 From Subsumption Architectures to Behavior-Based Robotics</b>	423
Subsumption Architectures: The Example of Allen	424
Behavior-Based Robotics: TOTO	427
Multiagent Programming: The Nerd Herd	430
<b>17 Looking Ahead: Challenges and Opportunities</b>	437
<b>17.1 Exploring the Connectivity of the Brain: The Human Connectome Project and Beyond</b>	438
<b>17.2 Understanding What the Brain Is Doing When It Appears Not to Be Doing Anything</b>	439
<b>17.3 Neural Prosthetics</b>	440
<b>17.4 Cognitive Science and the Law</b>	441
<b>17.5 Autonomous Vehicles: Combining Deep Learning and Intuitive Knowledge</b>	442
<i>Glossary</i>	444
<i>Bibliography</i>	454
<i>Index for Cognitive Science (3rd edition)</i>	478



## BOXES

- 2.1** A Conversation with ELIZA 39
- 3.1** What Does Each Lobe Do? 69
- 3.2** Brain Vocabulary 72
- 4.1** Defining Sentences in Propositional Logic 102
- 6.1** Basins of Attraction in State Space 157
- 7.1** Basic of the Probability Calculus 174
- 7.2** Deriving Bayes's Rule 177
- 8.1** The Prisoner's Dilemma 215
- 15.1** A Typical Semantic Priming Experiment 384



## FIGURES

- 0.1** Connections among the cognitive sciences, as depicted in the Sloan Foundation's 1978 report. 4
- 0.2** Some of the principal branches of scientific psychology. 7
- 0.3** Levels of organization and levels of explanation in the nervous system. 8
- 0.4** The spatial and temporal resolution of different tools and techniques in neuroscience. 9
- 0.5** The “space” of contemporary cognitive science. 11
- 1.1** A rat in a Skinner box. 18
- 1.2** A fourteen-unit T-Alley maze. 19
- 1.3** A cross-maze. 20
- 1.4** Schematic representation of a Turing machine. 25
- 1.5** A sample phrase structure tree for the sentence “John has hit the ball.” 27
- 1.6** Donald Broadbent’s (1958) model of selective attention. 29
- 2.1** A question for SHRDLU about its virtual micro-world. 40
- 2.2** An algorithm for determining whether a given input is a sentence or not. 42
- 2.3** Algorithms for identifying noun phrases and verb phrases. 43
- 2.4** Procedure for applying the command CLEARTOP. 44
- 2.5** SHRDLU acting on the initial command to pick up a big red block. 45
- 2.6** SHRDLU completing instruction 3 in the dialog: “Find a block which is taller than the one you are holding and put it in the box.” 46
- 2.7** Examples of the three-dimensional figures used in Shepard and Metzler’s 1971 studies of mental rotation. 48
- 2.8** Results of Shepard and Metzler’s 1971 studies of mental rotation. 49
- 2.9** Examples of vertically and horizontally oriented objects that subjects were asked to visualize in Kosslyn’s 1973 scanning study. 52
- 2.10** Two images of a bucket. 56

- 2.11** Two examples of Marr's primal sketch, the first computational stage in his analysis of the early visual system. 57
- 2.12** An example of part of the 2.5D sketch. 58
- 2.13** An illustration of Marr's 3D sketch, showing how the individual components are constructed. 59
- 2.14** The place of the implementational level within Marr's overall theory. 60
- 2.15** An illustration of the hierarchical organization of the visual system, including which parts of the brain are likely responsible for processing different types of visual information. 61
- 3.1** The large-scale anatomy of the brain, showing the forebrain, the midbrain, and the hindbrain. 69
- 3.2** A vertical slice of the human brain, showing the cerebrum. 70
- 3.3** The division of the left cerebral hemisphere into lobes. 71
- 3.4** The primary visual pathway. 72
- 3.5** Image showing ventral (purple) and dorsal (green) pathways in the human visual system. 73
- 3.6** Design and results of Ungerleider and Mishkin's cross-lesion disconnection studies. 75
- 3.7** A generic three-layer connectionist network (also known as an artificial neural network). 78
- 3.8** Gorman and Sejnowski's mine/rock detector network. 80
- 3.9** Images showing the different areas of activation (as measured by blood flow) during the four different stages in Petersen et al.'s (1988) lexical access studies. 84
- 3.10** A flowchart relating areas of activation to different levels of lexical processing. 85
- 3.11** Neural area showing activity when subjects looked at pictures. 88
- 3.12** Neural areas where activation is correlated with levels of memory performance. 89
- 3.13** A microelectrode making an extracellular recording. 90
- 3.14** Simultaneous microelectrode and fMRI recordings from a cortical site showing the neural response to a pulse stimulus of 24 seconds. 92
- 4.1** A typical traveling salesperson problem. 104
- 4.2** The structure of Fodor's argument for the language of thought hypothesis. 114
- 4.3** Inside and outside the Russian room. 116
- 5.1** Schematic illustration of a typical neuron. 125

- 5.2** An artificial neuron. 126
- 5.3** Four different activation functions. 127
- 5.4** Illustration of a mapping function. 128
- 5.5** A single-layer network representing the Boolean function AND. 130
- 5.6** A single-layer network representing the Boolean function NOT. 131
- 5.7** The starting configuration for a single-layer network being trained to function as a NOT-gate through the perceptron convergence rule. 133
- 5.8** Graphical representations of the AND and XOR (exclusive-OR) functions, showing the linear separability of AND. 135
- 5.9** A multilayer network representing the XOR (exclusive-OR) function. 136
- 5.10** The computational operation performed by a unit in a connectionist model. 138
- 6.1** The trajectory through state space of an idealized swinging pendulum. 151
- 6.2** The state space of a swinging pendulum in a three-dimensional phase space. 152
- 6.3** Illustration of the Watt governor, together with a schematic representation of how it works. 155
- 6.4** An example of the computational approach to motor control. 160
- 6.5** The stage IV search task, which typically gives rise to the A-not-B-error in infants at around the age of 9 months. 162
- 6.6** An infant sitting for an A trial and standing for a B trial. 163
- 6.7** Applying the dynamical field model to the A-not-B error. 165
- 7.1** An illustration purporting to be of Thomas Bayes from a 1936 book on the history of life insurance. 172
- 7.2** A diagram showing the proportion of the probability space in which A is true; the proportion of the probability space in which B is true; and the intersection of A and B (which is the region where A and B are both true). 175
- 7.3** Four of the seven Gestalt principles of grouping, illustrated and explained. 180
- 7.4** Two examples of stimuli used to elicit binocular rivalry. 182
- 7.5** Two well-known ambiguous figures: Rubin's vase and the duck–rabbit illusion. 183
- 7.6** The principal pathways for saccade production. 192
- 7.7** Platt and Glimcher's probabilistic cued saccade task. 193
- 7.8** Activity of an LIP neuron during the probability experiment. 194

- 7.9** Platt and Glimcher's cued saccade experiment, with stimulus and response held constant and the quantity of reward varied. 195
- 7.10** Activity of an LIP neuron while a monkey makes his own choice compared to a behaviorally derived estimate of the value of the movement to the monkey. 197
- 8.1** The architecture of a simple reflex agent. 205
- 8.2** The architecture of a goal-based agent. 206
- 8.3** The architecture of a learning agent. 207
- 8.4** A version of the Wason selection task. 212
- 8.5** A version of Griggs and Cox's deontic selection task. 213
- 8.6** The evolutionary biologist W. D. Hamilton (1936–2000). 217
- 8.7** The ACT-R cognitive architecture. 221
- 9.1** Luria's (1970) diagram of the functional organization of the brain. 231
- 9.2** Map of the anatomy of the brain showing the four lobes and the Brodmann areas. 233
- 9.3** A connectivity matrix for the visual system of the macaque monkey. 235
- 9.4** An anatomical wiring diagram of the visual system of the macaque monkey. 236
- 9.5** The results of single-neuron recordings of a mirror neuron in area F5 of the macaque inferior frontal cortex. 238
- 9.6** Typical patterns of EEG waves, together with where/when they are typically found. 239
- 9.7a** Common experimental design for neurophysiological studies of attention. 243
- 9.7b** Example of the occipital ERPs recorded in a paradigm of this nature. 244
- 9.7c** Single-unit responses from area V4 in a similar paradigm. 245
- 9.7d** Single-unit responses from area V1 showing no effect of attention. 245
- 9.8** Frontoparietal cortical network during peripheral visual attention. 247
- 9.9** An illustration of a typical delayed saccade task. 248
- 9.10** Peripheral attention versus spatial working memory versus saccadic eye movement across studies. 250
- 10.1** The dual-route model of past tense learning in English proposed by Steven Pinker and Alan Prince. 269
- 10.2** Rumelhart and McClelland's model of past tense acquisition. 270
- 10.3** Performance data for Rumelhart and McClelland's model of past tense learning. 271

- 
- 10.4** The network developed by Plunkett and Marchman to model children's learning of the past tense. 272
  - 10.5** A comparison of the errors made by Adam, a child studied by the psychologist Gary Marcus, and the Plunkett–Marchman neural network model of tense learning. 273
  - 10.6** A hierarchical cluster of similarity judgments, with nodes corresponding to clusters of stimuli more similar on average to each other than to objects in the nearest cluster. 279
  - 11.1** Schematic representation of the habituation and test conditions in Baillargeon's drawbridge experiments. 288
  - 11.2** Schematic representation of an experiment used to test infants' understanding of object boundaries and sensitivity to what Spelke calls the principle of cohesion. 289
  - 11.3** Schematic representation of an experiment testing infants' understanding of the principle of contact. 290
  - 11.4** Schematic depiction of events that accord with, or violate, the continuity or solidity constraints. 291
  - 11.5** A series of inputs to the network as a barrier moves in front of a ball and then back to its original location. 295
  - 11.6** Recurrent network for learning to anticipate the future position of objects. 296
  - 11.7** A balance beam. 297
  - 11.8** The architecture of the McClelland and Jenkins network for the balance beam problem. 299
  - 12.1** A decision tree illustrating a mortgage expert system. 309
  - 12.2** The first node on the decision tree for the tennis problem. 313
  - 12.3** The complete decision tree generated by the ID3 algorithm. 313
  - 12.4** A sample completed questionnaire used as input to an ID3-based expert system for diagnosing diseases in soybean crops. 314
  - 12.5** Different ways of distinguishing two groups in a database of examples. 317
  - 12.6** An illustration of hierarchical visual processing. 320
  - 12.7** Illustration of how an autoencoder compresses and then decompresses a signal. 323
  - 12.8** A move in the Google DeepMind challenge between AlphaGo and Lee Sedol in 2016. 328
  - 13.1** An example of metarepresentation. 338
  - 13.2** The general outlines of Leslie's model of pretend play. 339
  - 13.3** Leslie's Decoupler model of pretense. 340

- 13.4** The task used by Baron-Cohen, Leslie, and Frith to test for children's understanding of false belief. 344
- 13.5** Illustration of the connection between pretend play and success on the false belief task. 346
- 13.6** Baron-Cohen's model of the mindreading system. 350
- 14.1** What goes on when one subject represents another's belief. 361
- 14.2** What goes on when one subject represents another's perception. 362
- 14.3** A schematic version of standard simulationism. 364
- 14.4** Schematic representation of brain regions associated with the attribution of mental states. 367
- 14.5** Schematic overview of the frontoparietal mirror neuron system (MNS) and its main visual input in the human brain. 372
- 15.1** An illustration of a typical priming experiment. 382
- 15.2** Examples of deficits found in patients with left spatial neglect (damage to the right hemisphere of the brain). 386
- 15.3** D.B.'s responses to pictures of animals presented in his blind field. 388
- 15.4** An illustration of the two houses presented to P.S. 389
- 15.5** In this experiment, subjects were asked either to "post" a card into a slot or to rotate another hand-held card to match the orientation of the slot. 391
- 15.6** In the Ebbinghaus illusion, two circles are illusorily seen as differently sized, depending on what surrounds them. 392
- 15.7** In the Norman and Shallice 1980 model, conscious processing is involved in the supervisory attentional regulation, by prefrontal cortices, of lower-level sensorimotor chains. 398
- 15.8** The neural substrates of the global workspace. 399
- 16.1** A map of SHAKEY's physical environment. 409
- 16.2** A photograph of SHAKEY the robot. 410
- 16.3** The organizing principles of biorobotics – a highly interdisciplinary enterprise. 419
- 16.4** The cricket's ears are on its front legs. 420
- 16.5** A robot fish called WANDA. 421
- 16.6** WANDA swimming upward. 422

- 16.7** Another example of morphological computation: the robot hand designed by Hiroshi Yokoi. 422
- 16.8** The Yokoi hand grasping two very different objects. 423
- 16.9** Rodney Brooks's robot Allen, his first subsumption architecture robot. 424
- 16.10** The layers of Allen's subsumption architecture. 425
- 16.11** The Nerd Herd, together with the pucks that they can pick up with their grippers. 430



## TABLES

- 2.1** A table illustrating the three different levels that Marr identified for explaining information-processing systems 55
- 4.1** Syntax and semantics in the predicate calculus 113
- 8.1** Comparing the symbolic and subsymbolic dimensions of knowledge representation in the hybrid ACT-R architecture 224
- 9.1** Comparing techniques for studying connectivity in the brain 241
- 10.1** The stages of past tense learning according to verb type 268
- 13.1** The three groups studied in Baron-Cohen, Leslie, and Frith (1985) 343
- 16.1** SHAKEY'S five levels 411
- 16.2** How SHAKEY represents its own state 412
- 16.3** SHAKEY's intermediate-level routines 413
- 16.4** The five basis behaviors programmed into Matarić's Nerd Herd robots 431



## PREFACE



### About This Book

There are few things more fascinating than the human mind – and few things that are more difficult to understand. Cognitive science is the enterprise of trying to make sense of this most complex and baffling natural phenomenon.

The very things that make cognitive science so fascinating make it very difficult to study and to teach. Many different disciplines study the mind. Neuroscientists study the mind's biological machinery. Psychologists directly study mental processes, such as perception and decision-making. Computer scientists explore how those processes can be simulated and modeled in computers. Evolutionary biologists and anthropologists speculate about how the mind evolved. In fact, very few academic areas are not relevant to the study of the mind in some way. The job of cognitive science is to provide a framework for bringing all these different perspectives together.

The enormous range of information out there about the mind can be overwhelming, both for students and for instructors. Different textbooks have approached this challenge in different ways.

Some textbooks have concentrated on being as comprehensive as possible, with a chapter covering key ideas in each of the relevant disciplines – a chapter on psychology, a chapter on neuroscience, and so on. These books are often written by committee – with each chapter written by an expert in the relevant field. These books can be very valuable, but they really give an introduction to the cognitive sciences (in the plural) rather than to cognitive science as an interdisciplinary enterprise.

Other textbook writers take a much more selective approach, introducing cognitive science from the perspective of the disciplines that they know best – from the perspective of philosophy, for example, or of computer science. Again, I have learned much from these books, and they can be very helpful. But I am convinced that students and instructors need something more general.

This book aims for a balance between these two extremes. Cognitive science has its own problems and its own theories. The book is organized around these. They are all ways of working out the fundamental idea at the heart of cognitive science – which is that the mind is an information processor. What makes cognitive science so rich is that this single basic idea can be (and has been) worked out in many different ways. In presenting these different models of the mind as an information processor, I have tried to select as

wide a range of examples as possible to give students a sense of cognitive science's breadth and range.



## About the Third Edition

*Cognitive Science: An Introduction to the Science of the Mind* has been very significantly revised for the third edition. These changes have been made for two reasons. First, I wanted to make the book more accessible to students in the first and second years of their studies. To achieve that goal, I have made changes to both content and style as well as to the organization of the book. Second, I wanted to make sure that the new edition reflects the most exciting new developments in cognitive science, some of which were barely discernible back in 2010, when the first edition was published.

Previous editions of this book were organized around what I termed the *integration challenge*. This is the challenge of providing a unified theoretical framework for studying cognition that brings together the different disciplines that study the mind. The third edition no longer uses the integration challenge as an organizing principle. The additional layer of complexity is useful for many purposes, but not, I now think, for pedagogical ones. As a result, I have cut the two chapters that were devoted to the integration challenge in the first and second editions and simplified the presentation in later chapters. In particular, I no longer employ a two-way division into symbolic and nonsymbolic models of information processing. I have added an introduction that explains in some more general terms some of the issues and problems previously discussed under the label “integration challenge.”

I have used the space freed up by reorganization to expand coverage of more up-to-date areas elsewhere in the book. This includes a new chapter on Bayesian approaches to the mind. This chapter covers both the idea that cognition can be understood in terms of Bayesian hypothesis testing and error minimization and experimental studies in neuroeconomics of how probabilities and values seem to be calculated in a broadly Bayesian manner in the primate nervous system. In addition, I have updated the discussion of machine learning in what is now Chapter 12, eliminating some by now dated examples and replacing them with more topical discussion of *deep learning* algorithms.

To help instructors and students, I have divided some of the longer chapters from the second edition into two. Dynamical systems theory has its own chapter (Chapter 6), while situated cognition and robotics are now in Chapter 16. The lengthy discussion of mindreading in the second edition is now spread over two chapters: “Exploring Mindreading” (Chapter 13) and “Mindreading: Advanced Topics” (Chapter 14).



## How the Book Is Organized

This book is organized into three parts.



## Part I: Historical Landmarks

Cognitive science has evolved considerably in its short life. Priorities have changed as new methods have emerged – and some fundamental theoretical assumptions have changed with them. The three chapters in Part I introduce students to some of the highlights in the history of cognitive science. Each chapter is organized around key discoveries and/or theoretical advances.



## Part II: Models and Tools

Part II sets out the main models and tools that cognitive scientists can bring to bear to understand cognition and the mind.

The first model, discussed in Chapter 4, is associated with the physical symbol system hypothesis originally developed by the computer scientists Allen Newell and Herbert Simon. According to the physical symbol system hypothesis, all information processing involves the manipulation of physical structures that function as symbols. For the first decades of cognitive science, the physical symbol systems hypothesis was, as Jerry Fodor famously put it, the “only game in town.” In the 1980s and 1990s, connectionist and neural network modelers developed an alternative, derived from models of artificial neurons in computational neuroscience and connectionist artificial intelligence. Chapter 5 explores the motivation for this approach and introduces some of the key concepts.

Another set of models and tools that can be used to study the mind derives from dynamical systems theory and is introduced and discussed in Chapter 6. Bayesian approaches to modeling the mind have also gained currency. As explained in Chapter 7, these approaches treat the mind as a predictive, hypothesis-testing machine and have been used both to study the mind as a whole and to model the activity of individual brain areas and populations of neurons.

One of the key ideas of cognitive science is that the mind is modular (that some, or all, information processing is carried out by specialized modules). Chapter 8 explores different ways of developing this basic idea, including the radical claim, proposed by evolutionary psychologists, that the mind is simply a collection of specialized modules, with no non-specialized processing at all. Theoretical discussions of modularity are complemented by experimental techniques for studying the organization of the mind. Chapter 9 surveys the cognitive scientist’s tool kit in this regard, focusing in particular on different types of brain mapping.



## Part III: Applications

The seven chapters in this part are more applied than those in Part II. They explore different ways in which the models and tools introduced in Part II can be used to give accounts of particular cognitive phenomena.

Chapter 10 considers language learning. Many cognitive scientists have thought that language learning is a process of learning explicit rules with a significant innate component. But we explore both neural network and Bayesian models that illustrate alternative ways of thinking about how children can learn languages. In Chapter 11 we turn to models of how children learn about the basic structure of the physical world (how they acquire what is often called a *folk physics*). Here, too, we see the power of neural network models.

One of the most significant recent developments of neural network models has been the explosive growth of deep learning algorithms, which have made possible impressive advances in areas long thought to be major challenges for artificial intelligence, such as image recognition, machine translation, and games of strategy, such as Go. These are covered in Chapter 12.

Chapters 13 and 14 illustrate how theoretical, methodological, and experimental issues can come together. They work through an issue that has received much attention in contemporary cognitive science – the issue of whether there is a dedicated cognitive system response for our understanding of other people (the so-called mindreading system). Chapter 13 presents some of the basic issues and developments, while more advanced topics are introduced in Chapter 14.

In Chapter 15 we look at recent developments in the cognitive science of consciousness – a fast-moving and exciting area that raises fundamental questions about possible limits to what can be understood through the tools and techniques of cognitive science. And then finally, in Chapter 16, we explore the situated cognition movement and related developments in robotics, particularly behavior-based robotics and biologically inspired robotics.



## Using This Book in Courses

This book has been designed to serve as a self-contained text for a single-semester (12–15 weeks) introductory course on cognitive science. Students taking this course may have taken introductory courses in psychology and/or philosophy, but no particular prerequisites are assumed. All the necessary background is provided for a course at the freshman or sophomore level (first or second year). The book could also be used for a more advanced introductory course at the junior or senior level (third or fourth year). In this case, the instructor would most likely want to supplement the book with additional readings. There are suggestions on the instructor website (see below).



## Text Features

I have tried to make this book as user-friendly as possible. Key text features include the following:

■ **Chapter overviews.** Each chapter begins with an overview to orient the reader.

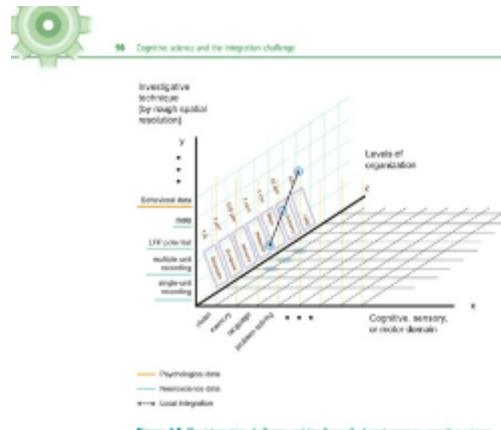


Figure 4.5 The integration challenge and the “space” of contemporary cognitive science.

In any event, whether the integration challenge is ultimately soluble or not, it is very clear that, as things stand, we are nowhere near solving it. Even the most ambitious theories and studies that have been carried out by cognitive scientists set out to cover only a tiny region of the space across which cognitive science ranges. Marr's theory of vision is one of the most ambitious undertakings of cognitive science, and Marr's level hypothesis is often cited as a textbook example of how cognitive science can span different levels of explanation. But the scope of Marr's theory is really just a very small part of vision. Marr's theory of vision is ultimately a theory of *visual* processing. It has nothing to say about object recognition, or about identification, nor about how vision is integrated with other sensory modalities or how memory and what is stored in memory. So, Marr's theory of vision can only say something about what you might think of as the *maths* of cognitive science – alternatively, it can explain only a very restricted historical slice of cognitive science. Moving to the x-axis, Marr had relatively little to say about what he called the ‘implementational level’. And, in fact, as we shall see in the next chapter (Section 5.2), the very idea that there is a single implementational level is deeply flawed.

**Local integration 1: Evolutionary psychology and the psychology of reasoning**

Cognitive psychologists have paid close attention to human problem-solving. We have already seen an example of this in the experiments on mental imagery and mental rotation. The issue here was how people solve problems that are framed in linguistic terms – problems involving the congruence of two figures, for example, from visual attention has been paid to problems that are logically stated, such as problems where subjects have to determine how likely it is that certain propositions are true, or whether one proposition follows from or is entailed by another. These problems are all reasoning problems, and psychologists have studied them with the aim of uncovering the principles of cognition.

A natural hypothesis in this area (particularly from those who have set themselves up on logic and probability theory) is that human reasoning is governed by the basic principles of logic and probability theory. People exploit the basic principles of logic when they are trying to solve problems that have a definite “yes/no” answer based on logical relations between propositions, and they use the principles of probability theory when the problem is to work out how likely some event is to happen. This may seem too obvious to be worth stating. How could one use anything but logic to solve logic problems? And how could we use anything but probability theory to solve probability problems?

Actually, however, the hypothesis is far from obviously true. Logic and probability theory are branches of mathematics, not of psychology. They study abstract mathematical relations. Those abstract mathematical relations determine the correct solution to particular problems. But logic and probability theory have nothing to say about how we actually go about solving those problems. In order to work out the reasoning principles that we actually use, psychologists have devised experiments to work out the sorts of problems that we are good at and the sorts of problems that we are bad at.

Before going on to look at some of those experiments we need to make explicit an important feature of both logic and probability theory. The basic laws of logic and principles of probability theory are universal. Logical relations hold between sentences irrespective of what those sentences actually say. We might, for example, make the following inference: “that’s the cathedral, then the library must be over there, but it’s not, so that can’t be the cathedral!” The logical rule here is known as *exclusis vellet*. This is the rule stating that if a conditional (*If A then B*) and the negation of the consequent of that conditional (not-B) jointly entail the negation of the antecedent of that conditional (not-A).

In our example the sentence “that’s the cathedral” takes the place of *A*, the antecedent of the conditional, and “the library must be over there” takes the place of *B*, the consequent of the conditional. What is distinctive about this sort of inference is that it makes no difference what sentence one puts in place of *A* and *B* in the standard terminology: this inferential transition is *domain-general*. Whatever one puts in place of *A* and *B* the

■ **Exercises.** These have been inserted at various points within each chapter. They are placed in the flow of the text to encourage the reader to take a break from reading and engage with the material. They are typically straightforward, but for a few, I have placed suggested solutions on the instructor website (see below).



of representational primitives and possible parameters of variation. Once again, it is easy to see why ‘informational encapsulation will secure computational tractability’! An informationally encapsulated module will have only a limited range of inputs to which to work.

In contrast, non-modular processing runs very quickly via versions of the so-called ‘place problem’. This is the problem, particularly pressing for those developing expert systems in AI and designing robots, of building into a system rules that will correctly identify what information and which information should be passed on in a given situation. The problem is identifying what sort of information is relevant and hence needs to be taken into account. Daniel Dennett's classic article on the subject opens with the following amusing and instructive tale:

Once upon a time there was a robot, named R1 by its creators. Its only task was to find for itself, one day a dragon arranged for it to learn that its spare battery, its precious energy supply, was locked in a room with a sturdily bolted door. R1 located the room, and the key to the door, and formulated a plan to rescue its battery. There was a wagon in the room, and the battery was on the wagon, and R1 hypothesized that a certain action which it called PULLOUT (Dragon, Room, 0) would result in the battery being removed from the room. Straightaway it acted, and did succeed in getting the battery out of the room before the bolts went off. Unfortunately, however, the bolts were on the wagon, and there was no way to roll the wagon in the room, but R1 didn't realize that driving the wagon along with the battery, R1, had started driving the wagon along with the battery.

Back to the drawing board. “The solution is obvious,” said the dragon. “Our new robot must be made to recognize not just the intended implications of its acts, but also the implications about their side-effects, by deducing those implications from the descriptions it uses in formulating its plans.” They called their new model, the robot-dragon, R2D2. They placed R2D2 in much the same predicament that R1 had succumbed to, and as the lit lit open the idea of PULLOUT (Dragon, Room, 0) it began, as designed, to consider the implications of such a course of action. It had just finished deducing that pulling the wagon out of the room would change the colour of the room's walls, and was on the cusp of a proof of the further implication that pulling the wagon out would cause it's wheels to hurt more than when there were wheels on the wagon – when the bath exploded.

Back to the drawing board. “We must teach it the difference between relevant implications and irrelevant implications,” said the dragon, “and, with it, ignore the irrelevant ones.” So they developed a method of ruling implications as either relevant or irrelevant to the project at hand, and installed the method in their new model, the relevance-detector, or R2D2 for short. When they subjected R2D2 to the test that had so ungraciously selected its ancestors for extinction, they were surprised to see nothing. Besides, outside the room containing the bickering hosts, the native life of its evolution stalked on with the gait of old thought, as Shakespeare (and more recently Fedor) has aptly put it: “Do something!” they prodded it. “Lanc,” it retorted. “I am hardly growing feathers thousands of implications I have determined to be irrelevant. Just as soon

as I find an irrelevant implication, I put it on the list of those I must ignore, and...” the bickering went on.

The greater the range of potentially relevant information, the more intractable this problem will be. This means that the tractability of the basic problem is in inverse proportion to the degree of information encapsulation. The more informationally encapsulated an informational system is, the less significant the frame problem will be. In the case of strictly modular systems, the frame problem will be negligible. In contrast, the less informationally encapsulated a system is, the more significant the frame problem will be. For non-modular systems, the frame problem has proven very hard indeed to tackle.

**P** Exercise 5.7 Explain in your own words what the frame problem is, without reference to the other examples. Differentiate the three approaches to the problem that Dennett sketches in the passage (again without reference to the soler example) and explain the difficulty with each of them.

For these two reasons, then, it looks very much as if the type of low-level, algorithmic analysis proposed by Marr works best for cognitive systems that are specialized, domain-specific, and informationally encapsulated – that is, for modular systems. And even if it could be extended to systems that are non-modular, Marr's approach would still not be applicable to the mind as a whole. Whether or not it is possible to provide a functional specification suitable to algorithmic formulation for high-level cognitive systems, it is hard to imagine what a functional specification would look like for the mind as a whole. But in the last analysis, an understanding of the mind as a whole is what a solution to the integration challenge is ultimately aiming at.

**5.3 Models of mental architecture**

In this section we explore an alternative approach to the integration challenge – one that provides a much better fit with what is actually going on in contemporary cognitive science than either of the two global approaches we have been considering. The infer-theoretic reduction approach and the tiled-world hypothesis both tackle the integration problem head-on. This section, however, tries very seriously the idea that cognitive science spans several levels of explanation and they each propose a different model for connecting activity at those different levels. The approach we will be exploring in this section tackles the problem from a different direction. It turns off from a basic assumption common to all the cognitive sciences and that shows how different ways of interpreting that basic assumption generate different models of the mind as a whole. These different models of the mind as a whole are what I am calling different *worlds*: *worlds* of mental architecture is a way of specifying the different components and levels of cognitive science.

- **Boxes.** Boxes have been included to provide further information about the theories and research discussed in the text. Readers are encouraged to work through these, but the material is not essential to the flow of the text.

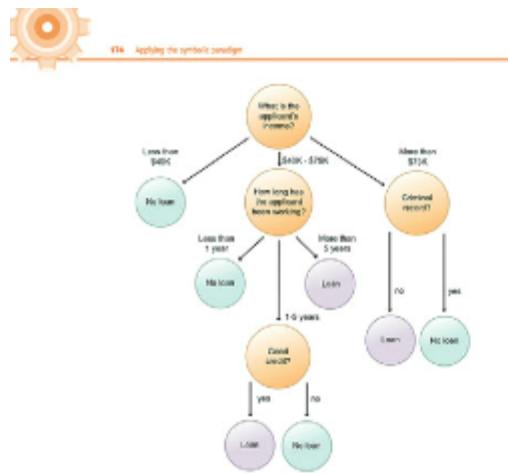


Figure 7.1 A decision tree illustrating a mortgage expert system. (from Friedenberg and Silverman 2006)

leads to a unique outcome (which computer scientists call a *terminal leaf* or *node*). But how are we supposed to get to these questions? How does the decision tree get designed, and how?

One simple way of doing it would be to ask a series of mortgage loan officers to sit down and come up with a decision tree that would map as best as possible to the practices at their bank. This could then be used as the basis for writing a program in a suitable programming language. This would be fine, and it is no doubt how many expert systems programs are actually written specifically in the mortgage area, but from the perspective of AI this would not be very interesting. It would be an expert system only in a very derivative sense. The real expert system would be the *space of mortgage*

### 7.1 Expert systems and machine learning

115

loan professionals. Much more interesting would be a program that was capable of producing its own decision tree – a program capable of imposing its own structure upon the problem and working out what would count as a solution. How would this work?

Here is one possible way of characterizing the problem. Suppose that we have a huge together with all the loans decisions that the bank has taken over a long period of time information about the applicants – their income, work history, credit rating, and so on. If we can find a way of representing the bank's past decisions in the form of a decision tree, so that each branch of the tree ends either in the loan being given or the loan being declined, then we can use that decision tree to “process” new applications.

#### ■ Machine learning and the physical symbol system hypothesis

We can put the same point the other way around. The decision tree in Figure 7.1 is a tool for analyzing new loan applications. The information that any applicant provides in response to the question that the tree poses will allow the applicant to choose one of the branches and the applicant will end up with their application either being approved or turned down. So the challenge for the expert system is to come up with a decision tree like that in Figure 7.1 from a situation of previous loan applicants, their personal information, and the decision that was eventually made.

This is a classic example of the type of problem tackled in the branch of AI known as *machine learning* (a subfield in expert systems research). The challenge is to produce an algorithm that will organize a complex database in terms of some attribute we are particularly interested in such as an applicant's loan worthiness, in the example we are considering. The organization takes the form of a decision tree, which will determine whether or not the attribute holds in a given case, whether or not the applicant is loan-worthy.

In the case of the mortgage loan decision tree the target attribute is labeled as *Loan*. All the branches of the decision tree must end in terminal nodes that have a value for the target attribute, i.e. they must say *Yes* or *No*. The decision tree is constructed by classifying the database in terms of other features such as *Good credit* for *Zero more than \$75k*. Once the decision tree has been constructed, it can then be used to decide whether some new instance (some new mortgage applicant) has the target attribute or not (i.e. is approved for the loan or not).

In the next section we will look in some detail at how an influential machine learning algorithm works, but first let me make explicit the connection with the physical symbol system hypothesis. As we saw on section 6.1, the physical symbol system hypothesis involves four basic claims:

<sup>1</sup> Symbols are physical patterns.

<sup>2</sup> Symbols can be combined to form complex symbol structures.

- **Summaries, checklists, and further reading.** These can be found at the end of each chapter. The summary provides a short overview of the chapter. The checklist allows students to review the key points of the chapter and also serves as a reference point for instructors. Suggestions of additional books and articles are provided to guide students' further reading on the topics covered in the chapter.



**CHAPTER THIRTEEN**

## New horizons: Dynamical systems and situated cognition

**OVERVIEW**

408

**13.1** Cognitive science and dynamical systems 409  
 What are dynamical systems? 409  
 The cybernetic systems approach 409  
 Cognitive systems 410  
 Representational systems 410

**13.2** Applying dynamical systems: two examples from child development 412  
 Two ways of thinking about motor control 412  
 By individual themes and then as an integrated view 412  
 Assessing the dynamical systems approach 413

**13.3** Situated cognition and situated interaction 419  
 The situated approach building a situated model 419  
 Situated cognition and knowledge representation 420  
 Representations 421  
 Blindsight 421  
 The philosophy of situated cognition 424

**13.4** From simulation to embodiment to behavior-based robotics 426  
 Behaviorism and behaviorism: The example of Alvin 421  
 Behavior-based robotics 422  
 Multi-agent programming: The Bird Herd 426

**Overview**

Throughout this book we have been working through some of the basic consequences of a single principle: this is the principle that cognition is information processing. It is in many ways the most important theoretical assumption of cognitive science. The main activities in Part I were to work out what this principle of information processing means in the context of the information-processing model of cognition in the middle of the twentieth century. In Part II we looked at different ways of thinking about information processing—the physical symbol system hypothesis and the neural networks model. Despite their very significant differences, the physical symbol system and neural networks approaches share a fundamental commitment to the idea that

403



## Course Website

A course website accompanies the book. It can be found at [www.cambridge.org/bermudez3](http://www.cambridge.org/bermudez3). This website contains

- a bank of test questions
- PowerPoint slides for each chapter, organized by section
- electronic versions of figures from the text
- review questions for each chapter that students can use to check their understanding and to review the material
- sample syllabi for courses of different lengths and different levels
- links to useful learning resources, videos, and experimental demonstrations
- links to online versions of relevant papers and online discussions for each chapter

Instructors can access a password-protected section of the website. This contains

- suggested solutions for the more challenging exercises and problems

The website is a work in progress. Students and instructors are welcome to contact me with suggestions, revisions, and comments. Contact details are on the website.

The screenshot shows the Academic Cambridge website interface. At the top, there's a navigation bar with 'Cart (0)' and 'Wishlist' buttons. Below that is a search bar with a magnifying glass icon and a checkbox for 'Include historic titles'. The main header features the word 'Academic' in large letters, with the tagline 'The future of publishing since 1584' underneath. A sub-header indicates the user is 'Welcome, Charles' and provides a 'Sign Out' link. The main menu includes 'Subjects', 'Textbooks', 'Reference', 'Authors', 'News', 'Conferences', 'Blogs', and 'Contact Us'. Below the menu, a breadcrumb trail shows the path: Home / Academic / Textbooks / Bermudez: Cognitive Science, 2nd edition. To the right of the trail are social media sharing icons. The central content area displays the book cover for 'José Luis Bermúdez Cognitive Science Second Edition'. The cover features a blue and green gear-themed design. To the left of the book cover is a sidebar with links for 'Home', 'For students' (Experimental demonstrations, Videos, Testbanks, Useful links), and 'For instructors' (PPT slides, Class discussion topics, Sample syllabi, Figures). A welcome message for the resources site is present, along with a brief description of the book's purpose and its second edition features. The footer contains a social media sharing bar.



## ACKNOWLEDGMENTS FOR THE FIRST EDITION

Many friends and colleagues associated with the Philosophy–Neuroscience–Psychology program at Washington University in St. Louis have commented on sections of this book. I would particularly like to thank Maurizio Corbetta, Frederick Eberhardt, David Kaplan, Clare Palmer, Gualtiero Piccinnini, Marc Raichle, Philip Robbins, David Van Essen, and Jeff Zacks. Josef Perner kindly read a draft of Chapter 12.

I have benefited from the comments of many referees while working on this project. Most remain anonymous, but some have revealed their identity. My thanks to Kirsten Andrews, Gary Bradshaw, Rob Goldstone, Paul Humphreys, and Michael Spivey.

Drafts of this textbook have been used four times to teach PNP 200 Introduction to Cognitive Science here at Washington University in St. Louis – twice by me and once each by David Kaplan and Jake Beck. Feedback from students both inside and outside the classroom was extremely useful. I hope that other instructors who use this text have equally motivated and enthusiastic classes. I would like to record my thanks to the teaching assistants who have worked with me on this course: Juan Montaña, Tim Oakberg, Adam Shriver, and Isaac Wiegman. And also to Kimberly Mount, the PNP administrative assistant, whose help with the figures and preparing the manuscript is greatly appreciated.

A number of students from my Spring 2009 PNP 200 class contributed to the glossary. It was a pleasure to work with Olivia Frosch, Katie Lewis, Juan Manfredi, Eric Potter, and Katie Sadow.

Work on this book has been made much easier by the efforts of the Psychology textbook team at Cambridge University Press – Raihanah Begum, Catherine Flack, Hetty Reid, Sarah Wightman, and Rachel Willsher (as well as to Andy Peart, who signed this book up but has since moved on). They have been very patient and very helpful. My thanks also to Anna Oxbury for her editing and to Liz Davey for coordinating the production process.



## ACKNOWLEDGMENTS FOR THE SECOND EDITION

I am very grateful to my colleagues in the Office of the Dean at Texas A&M University, particularly my administrative assistant Connie Davenport, for helping me to carve out time to work on the second edition of the textbook. T. J. Kasperbauer has been an excellent research assistant, providing numerous improvements to the text and supporting resources and helping me greatly with his deep knowledge of cognitive science. It has been a pleasure to work once again with Hetty Marx and Carrie Parkinson at Cambridge University Press. I particularly appreciate their work gathering feedback on the first edition. Thanks again to Anna Oxbury for her copyediting skills.



## ACKNOWLEDGMENTS FOR THE THIRD EDITION

Part of the work on the third edition was carried out during a period of leave in the 2018–19 academic year. I am grateful to the Department of Philosophy at Texas A&M University for allowing me to take a year from teaching and administrative duties. At CUP, I worked initially with Janka Romero and then subsequently with Stephen Acerra and Lisa Pinto, all of whom were very supportive and worked hard to get me helpful feedback on earlier editions. I am particularly grateful to the many reviewers who gave detailed comments and suggestions for improvement. Finally, I'd like to record my thanks to Dong An for her help with the online resources.





# Introduction

## The Challenge of Cognitive Science

### OVERVIEW 3

**0.1 Cognitive Science: An Interdisciplinary Endeavor 3**

**0.2 Levels of Explanation: The Contrast between Psychology and Neuroscience 5**

How Psychology Is Organized 6  
How Neuroscience Is Organized 7

**0.3 The Challenge of Cognitive Science 10**  
Three Dimensions of Variation 10  
The Space of Cognitive Science 10



## Overview

Cognitive science draws upon the tools and techniques of many different disciplines, including psychology, philosophy, linguistics, computer science, neuroscience, mathematical logic . . . It is a fundamentally *interdisciplinary activity*. This basic fact raises important and fundamental questions. What do all these disciplines have in common? How can they all come together to form a distinctive area of inquiry?

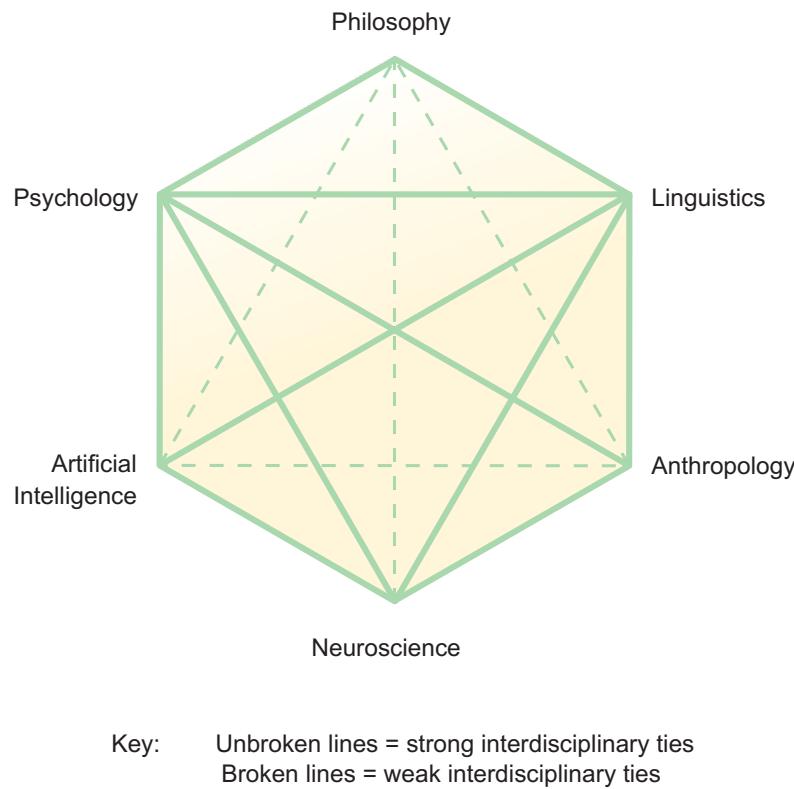
The aim of this introduction is to give you a sense of the scope and range of cognitive science, setting the framework for more detailed study in subsequent chapters. We will explore the idea that the different disciplines in cognitive science each study different levels of organization in the mind and the nervous system. In particular, we will see how the brain can be studied at many different levels, from the level of the molecule upward. The introduction ends with a description (and illustration) of what I call the space of cognitive science.



**0.1**

### Cognitive Science: An Interdisciplinary Endeavor

The hexagonal diagram in Figure 0.1 is one of the most famous images in cognitive science. It comes from the 1978 report on the state of the art in cognitive science commissioned by the Sloan Foundation and written by a group of leading scholars. The diagram is intended to illustrate the interdisciplinary nature of cognitive science. The lines on the diagram



**Figure 0.1** Connections among the cognitive sciences, as depicted in the Sloan Foundation's 1978 report. Unbroken lines indicate strong interdisciplinary links, while broken lines indicate weaker links. (Adapted from Gardner 1985)

indicate the academic disciplines that the authors saw as integral parts of cognitive science, together with the connections between disciplines particularly relevant to the study of mind and cognition.

For the authors of the Sloan report, cognitive science is an amalgamation of philosophy, psychology, linguistics, anthropology, neuroscience, and artificial intelligence. Each of the six disciplines brings with it different techniques, tools, and frameworks for thinking about the mind. Each of them studies the mind from different perspectives and at different levels. Whereas linguists, for example, develop abstract models of linguistic *competence* (the abstract structure of language), psychologists of language are interested in the mechanisms that make possible the *performance* of language users. Whereas neuroscientists study the details of how the brain works, computer scientists abstract away from those details to explore computer models and simulations of human cognitive abilities. Anthropologists are interested in the social dimensions of cognition, as well as how cognition varies across cultures. Philosophers, in contrast, are typically interested in very abstract models of how the mind is realized by the brain.

Some of the connections identified in the diagram were judged stronger than others. These are marked with a solid line. The weaker connections are marked with a broken line. At least one of the connections that was judged weak in 1978 has now become a



thriving subdiscipline in its own right. A group of philosophers impressed by the potential for fruitful dialog between philosophy and neuroscience have taken to calling themselves neurophilosophers, after the title of a very influential book by Patricia Churchland published in 1986.

Miller's own account of how the Sloan report was written is both disarming and telling. "The committee met once, in Kansas City. It quickly became apparent that everyone knew his own field and had heard of two or three interesting findings in other fields. After hours of discussion, experts in discipline X grew unwilling to make any judgments about discipline Y, and so forth. In the end, they did what they were competent to do: each summarized his or her own field and the editors – Samuel Jay Keyser, Edward Walker and myself – patched together a report" (Miller 2003: 143). This may be how reports get written, but it is not a very good model for an interdisciplinary enterprise such as cognitive science.

In fact, the hexagon as a whole is not a very good model for cognitive science. Even if we take seriously the lines that mark connections between the disciplines of cognitive science, the hexagon gives no sense of a unified intellectual enterprise. It gives no sense, that is, of something that is more than a composite of "traditional" disciplines such as philosophy and psychology. There are many different schools of philosophy and many different specializations within psychology, but there are certain things that bind together philosophers as a group and psychologists as a group, irrespective of their school and specialization. For philosophers (particularly in the so-called *analytic* tradition, the tradition most relevant to cognitive science), the unity of their discipline comes from certain problems that are standardly accepted as philosophical, together with a commitment to rigorous argument and analysis. The unity of psychology comes, in contrast, from a shared set of experimental techniques and paradigms. Is there anything that can provide a similar unity for cognitive science?

One of the main claims of this textbook is that cognitive science is indeed a unified enterprise. It has its own distinctive problems. Its own distinctive techniques, And its own distinctive explanatory frameworks. We will be studying all of these in this book. First, though, we need to get a better picture of the range and scope of the enterprise. In the rest of this introduction I'll use psychology and neuroscience as examples to give you a sense of the overall space of cognitive science.

## 0.2

## Levels of Explanation: The Contrast between Psychology and Neuroscience

Neuroscience occupies one pole of the Sloan report's hexagonal figure and it was not viewed as very central to cognitive science by the authors of the report. The report was written before the "turn to the brain" that we will look at in Chapter 3, and its focus reflected the contemporary focus on computer science, psychology, and linguistics as the core disciplines of cognitive science. Moreover, the authors of the report treated neuroscience as a unitary discipline, on a par with anthropology, psychology, and other more traditional academic disciplines. The explosion of research into what became known as cognitive neuroscience has since corrected both of these assumptions.

Most cognitive scientists place the study of the brain firmly at the heart of cognitive science. And it is becoming very clear that neuroscience is itself a massively interdisciplinary field.



## How Psychology Is Organized

One way of thinking about what distinguishes neuroscience from, say, psychology is through the idea of levels. I am talking here about what is sometimes called scientific psychology (psychology as it is taught and studied in university departments), as opposed, for example, to humanistic psychology, self-help psychology, and much of what is routinely classified as psychology in bookstores. But even narrowing it down like this, there are many different subfields of psychology.

A quick look at the courses on offer in any reputable psychology department will find courses in cognitive psychology, social psychology, abnormal psychology, personality psychology, psychology of language, and so on. It is normal for research psychologists to specialize in at most one or two of these fields. Nonetheless, most psychologists think that psychology is a single academic discipline. This is partly because there is a continuity of methodology across the different specializations and subfields. Students in psychology are typically required to take a course in research methods. Such courses cover basic principles of experimental design, hypothesis formation and testing, and data analysis that are common to all branches of psychology.

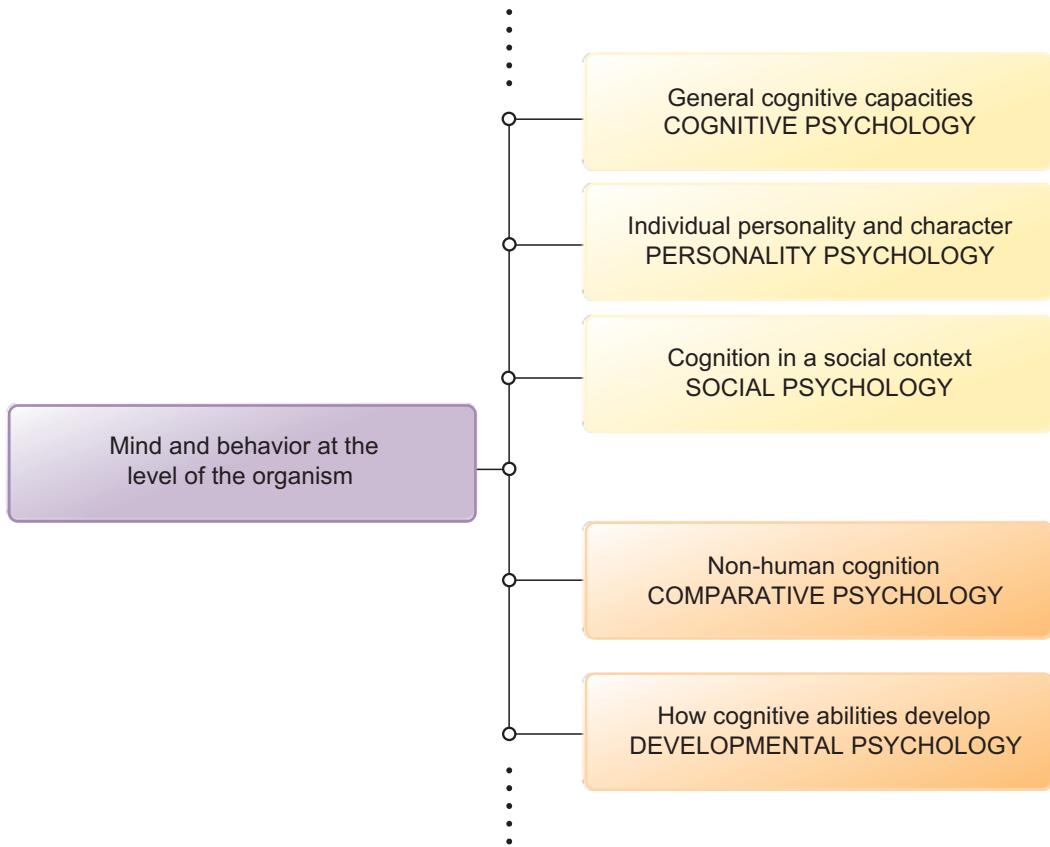
Equally important, however, is the fact that many of these branches of psychology operate at the same level. The data from which they begin are data about cognitive performance and behavior at the level of the whole organism (I am talking about the whole organism to make clear that these ideas extend to nonhuman organisms, as studied in comparative psychology).

The basic *explananda* (the things that are to be explained) in psychology are people's psychological capacities, which includes both cognitive and emotional capacities. The organization of psychology into different subfields is a function of the fact that there are many different types of cognitive and emotional capacities.

Within cognitive psychology, for example, what psychologists are trying to explain are the organism's capacities for perception, memory, attention, and so on. Controlled experiments and correlational studies are used to delimit and describe those capacities, so that psychologists know exactly what it is that needs to be explained.

Likewise, social psychologists study the capacities involved in social understanding and social interactions. They are interested, for example, in social influences on behavior, on how we respond to social cues, and on how our thoughts and feelings are influenced by the presence of others. Personality psychologists study the traits and patterns of behavior that go to make up what we think of as a person's character. And so on.

If we were to map out some of the principal subfields in scientific psychology it would look something like Figure 0.2. The diagram is intended to show that the different subbranches all study different aspects of mind and behavior at the level of the organism.



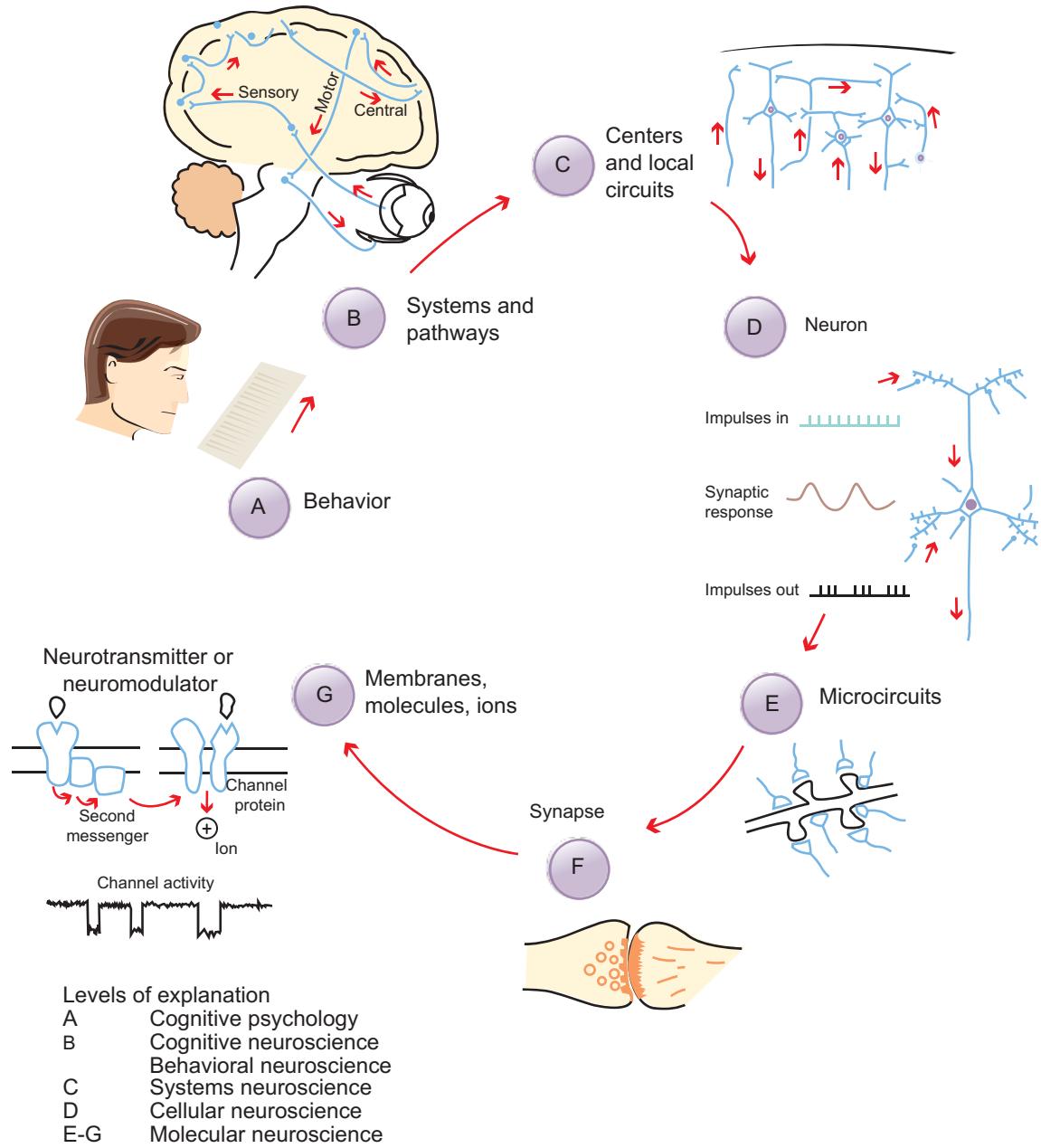
**Figure 0.2** Some of the principal branches of scientific psychology.

## How Neuroscience Is Organized

Things are very different in neuroscience. There are many branches of neuroscience, but they are not related in the same way. The organization of neuroscience into branches closely follows the different levels of organization in the brain and the central nervous system. These levels of organization are illustrated in Figure 0.3, drawn from Gordon Shepherd's 1994 textbook *Neurobiology*.

You may have come across references to areas in the brain such as the primary visual cortex or the hippocampus, for example. And you may have encountered talk of neural pathways connecting different areas in the brain. Located at levels A and B in Shepherd's diagram, these are the highest levels of neural organization, corresponding most closely to cognitive activities that we all perform. The primary visual cortex, for example, is responsible for coding the basic features of visual information coming from the retina. It is sensitive to orientation, motion, speed, direction, and so on. The hippocampus, in contrast, is thought to be responsible for key aspects of memory.

Activity at this top level of organization is the result of activity at lower levels of organization. In Shepherd's diagram this takes us to levels C and E – the level of centers, local circuits, and microcircuits. Somehow the collective activity of populations of neurons codes certain types of information about objects in a way that organizes and coordinates



**Figure 0.3** Levels of organization and levels of explanation in the nervous system. (Adapted from Shepherd 1994)

the information carried by individual neurons. These populations of neurons are the local circuits in Shepherd's diagram.

What happens in populations of neurons is ultimately determined by the behavior of individual neurons. But neurons are not the most basic level of organization in the nervous system. In order to understand how neurons work we need to understand how they communicate. This brings us to Shepherd's level F, because neurons communicate across synapses. Most synapses are chemical, but some are electrical. The chemical synapses work through the transmission of neurochemicals (*neurotransmitters*). These neurotransmitters are activated by

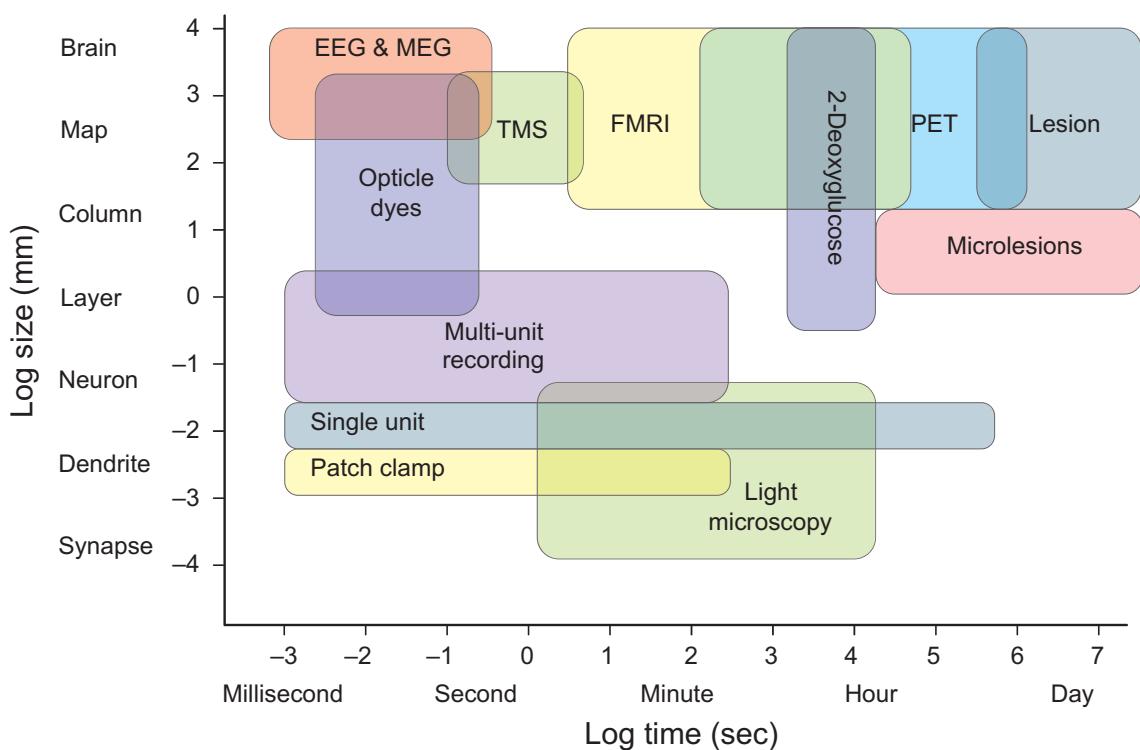


the arrival of an electrical signal (the *action potential*). The propagation of neurotransmitters works the way it does because of the molecular properties of the synaptic membrane – properties that are ultimately genetically determined. With this we arrive at level G in Shepherd's diagram.

The point of this whistle-stop tour through the levels of organization in the brain is that the subfields of neuroscience map very closely onto the different levels of organization in the brain. At the top level we have cognitive neuroscience and behavioral neuroscience, which study the large-scale organization of the brain circuits deployed in high-level cognitive activities. These operate at what in discussing the subfields of psychology I termed the level of the whole organism. Systems neuroscience, in contrast, investigates the functioning of neural systems, such as the visual system. The bridge between the activity of neural systems and the activity of individual neurons is one of the central topics in computational neuroscience, while cellular and molecular neuroscience deal with the fundamental biological properties of neurons.

Different branches of neuroscience (and cognitive science in general) employ tools appropriate to the level of organization at which they are studying the brain. These tools and techniques vary in what neuroscientists call their temporal and spatial resolution. That is, they vary in the scale on which they give precise measurements (spatial resolution) and the time intervals to which they are sensitive (temporal resolution).

Some of the important variations are depicted in Figure 0.4. We will explore the differences between these different tools and technologies in much more detail in later chapters (particularly Chapter 9).



**Figure 0.4** The spatial and temporal resolution of different tools and techniques in neuroscience. Time is on the x-axis and size is on the y-axis. (Adapted from Baars and Gage 2010)

## 0.3

### The Challenge of Cognitive Science

This section explores these basic ideas of levels of organization, levels of resolution, and levels of explanation further, to give a picture of what I call the space of cognitive science.

#### Three Dimensions of Variation

Cognitive science draws upon a large number of potentially relevant fields and subfields. Those fields and subfields differ from each other along three dimensions.

One dimension of variation is illustrated by the subfields of neuroscience. Neuroscience studies the brain at many different levels. These levels are organized into a vertical hierarchy that corresponds to the different levels of organization in the nervous system.

A second dimension of variation comes with the different techniques and tools that cognitive scientists can employ. As illustrated in Figure 0.4, these tools vary both in spatial and in temporal resolution. Some tools, such as PET and fMRI, give accurate measurements at the level of individual brain areas. Others, such as microelectrode recording, give accurate measurements at the level of individual neurons (or small populations of neurons).

The third dimension of variation is exemplified by the different subfields of psychology. Most of psychology operates at Shepherd's level A. The different areas of psychology set out to explore, map, describe, and explain are the cognitive abilities making possible the myriad things that human beings do and say.

#### The Space of Cognitive Science

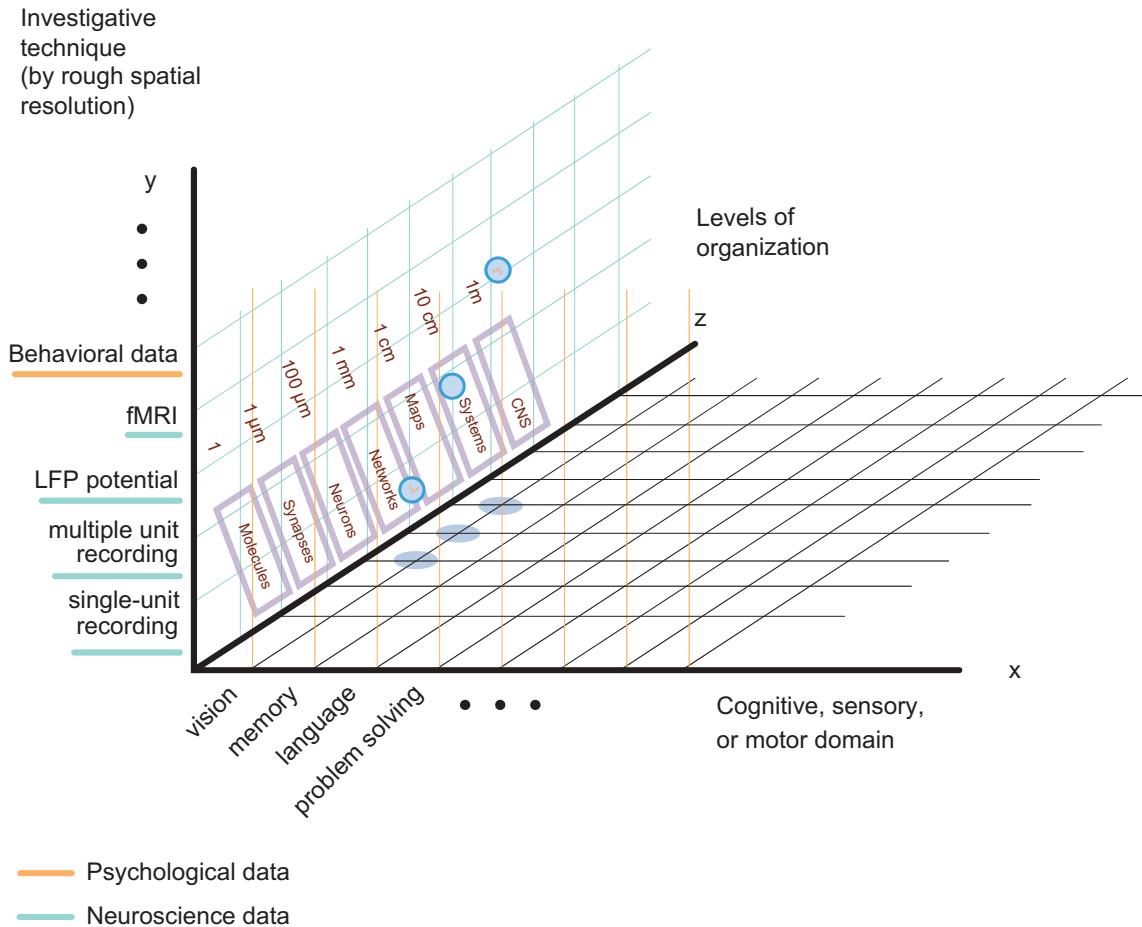
The different parts of cognitive science are distributed, therefore, across a three-dimensional space illustrated in Figure 0.5.

- The  $x$ -axis marks the different cognitive domains that are being studied
- The  $y$ -axis marks the different tools that might be employed (ordered roughly in terms of their degree of spatial resolution).
- The  $z$ -axis marks the different levels of organization at which cognition is studied.

This three-dimensional diagram is a more accurate representation of where cognitive science stands in the early years of the twenty-first century than the two-dimensional hexagon proposed by the authors of the Sloan report (although the hexagon may well have been an adequate picture of how things stood at the end of the 1970s).

A good way of thinking about cognitive science is as setting out to provide a unified account of cognition that draws upon and integrates the whole space. Cognitive science is more than just the sum of its parts. The aim of cognitive science as an intellectual enterprise is to provide a framework that makes explicit the common ground between all the different academic disciplines that study the mind and that shows how they are related to each other.

You can think of the analogy with physics. Many theoretical physicists think that the ultimate goal of physics is to provide a unified Theory of Everything. So too (on this way of thinking about cognitive science) is it the mission of cognitive science to provide a unified Theory of Cognition.



**Figure 0.5** The “space” of contemporary cognitive science.

Parts II and III will explore the principal theories of cognition in cognitive science, and see how they can be applied to explain different aspects of cognition. First, though, we turn to an overview of some of the key historical landmarks in the emergence and subsequent development of cognitive science. That will occupy the three chapters of Part I. These chapters should put flesh on the bones of the general picture sketched out in this introduction.

## Further Reading

Historical background on the Sloan report can be found in Gardner 1985 and Miller 2003 (available in the online resources). The report itself was never published. A very useful basic introduction to levels of organization and structure in the nervous system is chapter 2 of Churchland and Sejnowski 1992. For more detail, a classic neuroscience textbook is Kandel, Schwarz, and Jessell 2012. Stein and Stoodley 2006 and Purves et al. 2011 are alternatives. Craver 2007 discusses the interplay between different levels of explanation in the neuroscience of memory. Piccinini and Craver 2011 is a more general discussion; also see Bickle 2006 and Sullivan 2009.

PART I

# HISTORICAL LANDMARKS









## CHAPTER ONE

# The Prehistory of Cognitive Science

### OVERVIEW 15

- 1.1 The Reaction against Behaviorism in Psychology 16**  
Learning without Reinforcement: Tolman and Honzik, "‘Insight’ in Rats" (1930) 17  
Cognitive Maps in Rats? Tolman, Ritchie, and Kalish, "Studies in Spatial Learning" (1946) 20  
Plans and Complex Behaviors: Lashley, "The Problem of Serial Order in Behavior" (1951) 21
- 1.2 The Theory of Computation and the Idea of an Algorithm 22**  
Algorithms and Turing Machines: Turing, "On Computable Numbers, with an Application to the Decision Problem" (1936–7) 23

**1.3 Linguistics and the Formal Analysis of Language 25**

The Structure of Language: Chomsky’s *Syntactic Structures* (1957) 26

**1.4 Information-Processing Models in Psychology 28**

How Much Information Can We Handle? George Miller’s "The Magical Number Seven, Plus or Minus Two" (1956) 29

The Flow of Information: Donald Broadbent’s "The Role of Auditory Localization in Attention and Memory Span" (1954) and *Perception and Communication* (1958) 30

**1.5 Connections and Points of Contact 32**



## Overview

In the late 1970s cognitive science became an established part of the intellectual landscape. At that time an academic field crystallized around a basic set of problems, techniques, and theoretical assumptions. These problems, techniques, and theoretical assumptions came from many different disciplines and areas. Many of them had been around for a fairly long time. What was new was the idea of putting them together as a way of studying the mind.

Cognitive science is at heart an interdisciplinary endeavor. In interdisciplinary research great innovations come about simply because people see how to combine things that are already out there but have never been put together before. A good way to understand cognitive science is to try to think your way back to how things might have looked to its early pioneers. They were exploring a landscape in which certain regions were well mapped and well understood, but where

there were no standard ways of getting from one region to another. An important part of what they did was to show how these different regions could be connected in order to create an interdisciplinary science of the mind.

In this chapter we go back to the 1930s, 1940s, and 1950s – to explore the *prehistory* of cognitive science. We will be looking at some of the basic ideas and currents of thought that, in retrospect, we can see as feeding into what came to be known as cognitive science. As we shall see in more detail later on in this book, *the guiding idea of cognitive science is that mental operations involve processing information*, and hence that we can study how the mind works by studying how information is processed. This basic idea of the mind as an information processor has a number of very specific roots, in areas that seem on the face of it to have little in common. The prehistory of cognitive science involves parallel, and largely independent, developments in psychology, linguistics, and mathematical logic. We will be looking at four of these developments:

- The reaction against behaviorism in psychology (Section 1.1)
- The idea of algorithmic computation in mathematical logic (Section 1.2)
- The emergence of linguistics as the formal analysis of language (Section 1.3)
- The emergence of information-processing models in psychology (Section 1.4)

In concentrating on these four developments we will be passing over other important influences, such as neuroscience and neuropsychology. This is because until quite recently the direct study of the brain had a relatively minor role to play in cognitive science.

Almost all cognitive scientists are convinced that in some fundamental sense the mind just is the brain, so that everything that happens in the mind is happening in the brain. Few, if any, cognitive scientists are *dualists*, who think that the mind and the brain are two separate and distinct things. But for a long time in the history of cognitive science it was widely held that we are better off studying the mind by abstracting away from the details of what is going on in the brain. This changed only with the emergence in the 1970s and 1980s of new technologies for studying neural activity and of new ways of modeling cognitive abilities – as we will see in Chapter 3.

## 1.1

### The Reaction against Behaviorism in Psychology

Behaviorism was (and in some quarters still is) an influential movement in psychology. It takes many different forms, but they all share the basic assumption that psychologists should confine themselves to studying observable phenomena and measurable behavior. Behaviorists think that psychologists should avoid speculating about unobservable mental states, and instead focus on nonpsychological mechanisms linking particular stimuli with particular responses. These mechanisms are the product of conditioning. For examples of conditioning, think of Pavlov's dogs being conditioned to salivate at the sound of the bell, or the rewards/punishments that animal trainers use to encourage/discourage certain types of behavior.

For behaviorists, psychology is really the science of behavior. This approach to psychology leaves little room for cognitive science as the scientific study of cognition and the mind. Cognitive science could not even get started until behaviorism ceased to be the



dominant approach within psychology. Psychology's move from behaviorism was a lengthy and drawn-out process (and some would say that it has not yet been completed). We can appreciate some of the ideas that proved important for the later development of cognitive science by looking at three landmark papers. Each was an important statement of the idea that various types of behavior could not be explained in terms of stimulus-response mechanisms. Instead, psychologists need to think about organisms as storing and processing information about their environment, rather than as responding mechanically to reinforcers and stimuli. This idea of organisms as information processors is the single most fundamental idea of cognitive science.

## Learning without Reinforcement: Tolman and Honzik, "'Insight' in Rats" (1930)

Edward Tolman (1886–1959) was a behaviorist psychologist studying problem solving and learning in rats (among other things). As with most psychologists of the time, he started off with two standard behaviorist assumptions about learning. The first assumption is that all learning is the result of *conditioning*. The second assumption is that conditioning depends upon processes of *association* and *reinforcement*.

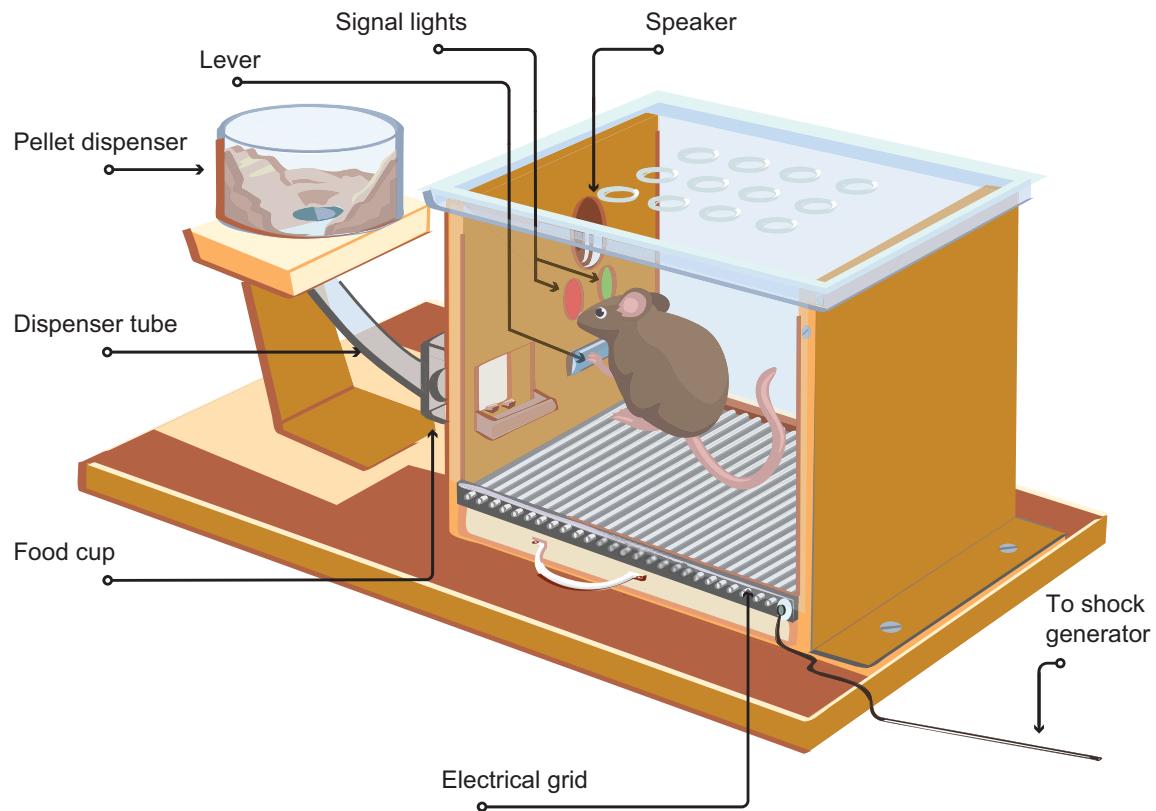
We can understand these two assumptions by thinking about a rat in what is known as a Skinner box, after the celebrated behaviorist B. F. Skinner. A typical Skinner box is illustrated in Figure 1.1. The rat receives a reward each time it behaves in a particular way (pressing a lever, for example, or pushing a button). The reward *reinforces* the behavior. This means that the association between the behavior and the reward is strengthened and the rat's performing the behavior again becomes more likely. The rat becomes *conditioned* to perform the behavior.

The basic idea of behaviorism is that all learning is either reinforcement learning of this general type, or the even simpler form of associative learning often called classical conditioning.

In classical conditioning what is strengthened is the association between a *conditioned stimulus* (such as the typically neutral sound of a bell ringing) and an *unconditioned stimulus* (such as the presentation of food). The unconditioned stimulus is *not* neutral for the organism and typically provokes a behavioral response, such as salivation. What happens during classical conditioning is that the strengthening of the association between conditioned stimulus and unconditioned stimulus eventually leads the organism to produce the unconditioned response to the conditioned stimulus alone, without the presence of the unconditioned stimulus. The most famous example of classical conditioning is Pavlov's dogs, who were conditioned to salivate to the sound of a bell by the simple technique of using the bell to signal the arrival of food.

So, it is a basic principle of behaviorism that all learning, whether by rats or by human beings, takes place through processes of reinforcement and conditioning. What the studies reported by Tolman and Honzik in 1930 seemed to show, however, is that this is not true even for rats.

Tolman and Honzik were interested in how rats learned to navigate mazes. They ran three groups of rats through a maze of the type illustrated in Figure 1.2. The first group



**Figure 1.1** A rat in a Skinner box. The rat has a response lever controlling the delivery of food, as well as devices allowing different types of stimuli to be produced.

received a reward each time they successfully ran the maze. The second group never received a reward. The third group was unrewarded for the first ten days and then began to be rewarded.

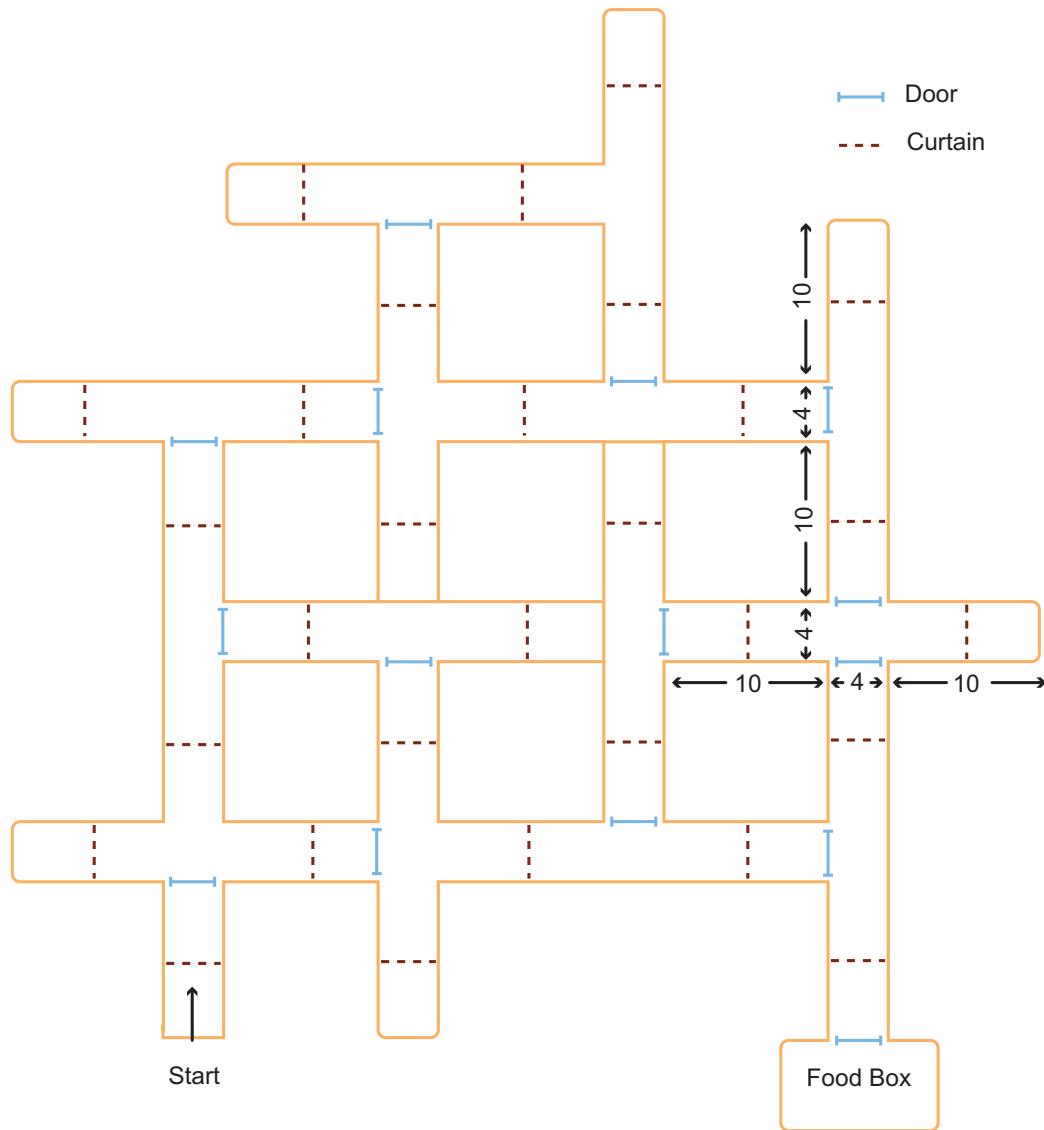
As behaviorism predicted, the rewarded rats quickly learned to run the maze, while both groups of unrewarded rats simply wandered around aimlessly. The striking fact, however, was that when the third group of rats started to receive rewards they learned to run the maze far more quickly than the first group had.

Tolman and Honzik argued that the rats must have been learning about the layout of the maze during the period when they were not being rewarded. This type of *latent learning* seemed to show that reinforcement was not necessary for learning, and that the rats must have been picking up and storing information about the layout of the maze when they were wandering around it, even though there was no reward and hence no reinforcement. They were later able to use this information to navigate the maze.



**Exercise 1.1** Explain in your own words why latent learning seems to be incompatible with the two basic assumptions of behaviorism.

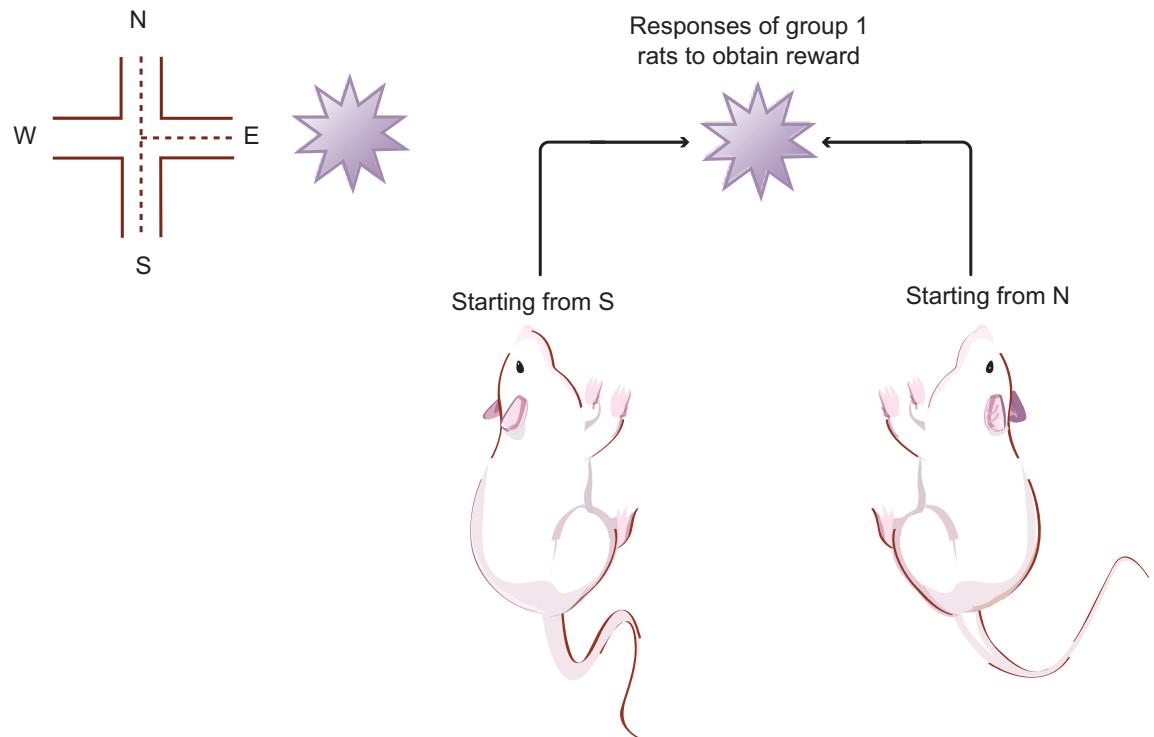
Suppose, then, that organisms are capable of latent learning – that they can store information for later use without any process of reinforcement. One important follow-up



**Figure 1.2** A fourteen-unit T-Alley maze (measurements in inches). Note the blocked passages and dead ends. (Adapted from Elliott 1928)

question is: What sort of information is being stored? In particular, are the rats storing information about the spatial layout of the maze? Or are they simply “remembering” the sequences of movements (responses) that they made while wandering around the maze? And so, when the rats in the latent-learning experiments start running the maze successfully, are they simply repeating their earlier sequences of movements, or are they using their “knowledge” of how the different parts of the maze fit together?

Tolman and his students and collaborators designed many experiments during the 1930s and 1940s to try to decide between *place learning* and *response learning* accounts of how rats learn to run a maze. Some of these experiments were reported in a famous article in 1946.



**Figure 1.3** A cross-maze, as used in Tolman, Ritchie, and Kalish (1946). The left-hand part of the figure illustrates the maze, with a star indicating the location of the food reward. The right-hand side illustrates how the group 1 rats had to make different sequences of movements in order to reach the reward, depending on where they started.

### Cognitive Maps in Rats? Tolman, Ritchie, and Kalish, "Studies in Spatial Learning" (1946)

One experiment used a cross-maze with four end points (North, South, East, West), like that illustrated in Figure 1.3. Rats were started at North and South on alternate trials. One group of rats was rewarded by food that was located at the same end point, say East. The relevant feature of the map for this group was that the same turning response would not invariably return them to the reward. To get from North to East the rat needed to make a left-hand turn, whereas a right-hand turn was required to get from South to East. For the second group the location of the food reward was shifted between East and West so that, whether they started at North or South, the same turning response was required to obtain the reward. A rat in the second group starting from North would find the reward at East, while the same rat starting from South would find the reward at West. Whether it started at North or South a left turn would always take it to the reward.

This simple experiment shows very clearly the distinction between place learning and response learning. Consider the first group of rats (those for which the food was always in the same place, although their starting-points differed). In order to learn to run the maze and obtain the reward they had to represent the reward as being at a particular place and



control their movements accordingly. If they merely repeated the same response they would only succeed in reaching the food reward on half of the trials. For the second group, though, repeating the same turning response would invariably bring them to the reward, irrespective of the starting point.

Tolman found that the first group of rats learned to run the maze much more quickly than the second group. From this he drew conclusions about the nature of animal learning in general, namely, that it was easier for animals to code spatial information in terms of places rather than in terms of particular sequences of movements.



### **Exercise 1.2 Explain in your own words why the experimental results seem to show that rats engage in place learning rather than response learning.**

Tolman took his place-learning experiments as evidence that animals form high-level representations of how their environment is laid out – what he called *cognitive maps*. Tolman's cognitive maps were one of the first proposals for explaining behavior in terms of *representations* (stored information about the environment).

Representations are one of the fundamental explanatory tools of cognitive science. Cognitive scientists regularly explain particular cognitive achievements (such as the navigational achievements of rats in mazes) by modeling how the organism is using representations of the environment. Throughout this book we will be looking at different ways of thinking about how representations code information about the environment, and about how those representations are manipulated and transformed as the organism negotiates and engages with its environment.

## **Plans and Complex Behaviors: Lashley, "The Problem of Serial Order in Behavior" (1951)**

At the same time as Tolman was casting doubt on standard behaviorist models of spatial navigation, the psychologist and physiologist Karl Lashley was thinking more generally about the problem of explaining complex behavior.

Much human and animal behavior has a very complex structure, involving highly organized sequences of movements. Stimulus-response behaviorists have limited resources for thinking about these complex behaviors. They have to view them as linked sequences of responses – as a sort of chain with each link determined by the link immediately preceding it. This is the basic idea behind response learning models of how rats run mazes.

The standard behaviorist view is that rats learn to chain together a series of movements that leads to the reward. Tolman showed that this is not the right way to think about what happens when rats learn to run mazes. Lashley made the far more general point that this seems to be completely the wrong way to think about many complex behaviors.

Think of the complicated set of movements involved in uttering a sentence of English, for example. Or playing a game of tennis. In neither of these cases is what happens at a particular moment solely determined by what has just happened – or prompted by what is going on in the environment and influencing the organism. What happens at any given

point in the sequence is often a function of what will happen later in the sequence, as well as of the overall goal of the behavior.

According to Lashley, we should think about many of these complex behaviors as products of prior planning and organization. The behaviors are organized hierarchically (rather than linearly). An overall plan (say, walking over to the table to pick up the glass) is implemented by simpler plans (the walking plan and the reaching plan), each of which can be broken down into simpler plans, and so on. Very little (if any) of this planning takes place at the conscious level.



### Exercise 1.3 Give your own example of a hierarchically organized behavior.

Lashley's essay contains the seeds of two ideas that have proved very important for cognitive science. The first is the idea that much of what we do is under the control of planning and information-processing mechanisms that operate below the threshold of awareness. This is the *hypothesis of subconscious information processing*. Even though we are often conscious of our high-level plans and goals (of what goes on at the top of the hierarchy), we tend not to be aware of the information processing that translates those plans and goals into actions. So, for example, you might consciously form an intention to pick up a glass of water. But carrying out the intention requires calculating very precisely the trajectory that your arm must take, as well as ensuring that your hand is open to the right degree to take hold of the glass. These calculations are carried out by information-processing systems operating far below the threshold of conscious awareness.

The second important idea is the *hypothesis of task analysis*. This is the idea that we can understand a complex task (and the cognitive system performing it) by breaking it down into a hierarchy of more basic subtasks (and associated subsystems). This hypothesis has proved a powerful tool for understanding many different aspects of mind and cognition. We can think about a particular cognitive system (say, the memory system) as carrying out a particular task – the task of allowing an organism to exploit previously acquired information. We can think about that task as involving a number of simpler, subtasks – say, the subtask of storing information and the subtask of retrieving information. Each of these subtasks can be carried out by even more simple sub-subtasks. We might distinguish the sub-subtask of storing information for the long term from the sub-subtask of storing information for the short term. And so on down the hierarchy.

## 1.2 The Theory of Computation and the Idea of an Algorithm

At the same time as Tolman, Lashley, and others were putting pressure on some of the basic principles of behaviorism, the theoretical foundations for one highly influential approach to cognitive science (and indeed for our present-day world of omnipresent computers and constant flows of digital information) were laid in the 1930s, in what was at the time a rather obscure and little-visited corner of mathematics.

In 1936–7 Alan Turing published an article in the *Proceedings of the London Mathematical Society* that introduced some of the basic ideas in the theory of computation. Computation



is what computers do and, according to many cognitive scientists, it is what minds do. What Turing gave us was a theoretical model that many have thought to capture the essence of computation. Turing's model (the so-called Turing machine) is one of the most important and influential ideas in cognitive science, even though it initially seems to have little to do with the human mind.

## Algorithms and Turing Machines: Turing, "On Computable Numbers, with an Application to the Decision Problem" (1936–7)

Turing, together with a number of mathematicians working in the foundations of mathematics, was grappling with the problem (known as the Halting Problem) of determining whether there is a purely mechanical procedure for working out whether certain basic mathematical problems have a solution.

Here is a way of understanding the Halting Problem. Think about it in terms of computer programs. Many computer programs are not defined for every possible input. They will give a solution for some inputs, the ones for which they are defined. But for other inputs, the ones for which they are not defined, they will just endlessly loop, looking for a solution that isn't there. From the point of view of a computer programmer, it is really important to be able to tell whether or not the computer program is defined for a given input – in order to be able to tell whether the program is simply taking a very long time to get to the solution, or whether it is in an endless loop.

This is what a solution to the Halting Problem would give – a way of telling, for a given computer program and a given input, whether the program is defined for that input. The solution has to work both ways. It has to give the answer "Yes" when the program is defined, and "No" when the program is not defined.

It is important to stress that Turing was looking for a purely mechanical solution to the Halting Problem. He was looking for something with the same basic features as the "recipes" that we all learn in high school for multiplying two numbers, or performing long division. These recipes are mechanical because they do not involve any insight. The recipes can be clearly stated in a finite set of instructions and following the instructions correctly always gives the right answer, even if you don't understand how or why.

Since the notion of a purely mechanical procedure is not itself a mathematical notion, the first step was to make it more precise. Turing did this by using the notion of an *algorithm*. An algorithm is a finite set of rules that are unambiguous and that can be applied systematically to an object or set of objects to transform it or them in definite and circumscribed ways. The instructions for programming a DVD recorder, for example, are intended to function algorithmically so that they can be followed blindly in a way that will transform the DVD recorder from being unprogrammed to being programmed to switch itself on and switch itself off at appropriate times. Of course, the instructions are not genuinely algorithmic since, as we all know, they are not idiot-proof.



### Exercise 1.4 Think of an example of a genuine algorithm, perhaps from elementary arithmetic or perhaps from everyday life.

One of Turing's great contributions was a bold hypothesis about how to define the notion of an algorithm within mathematics. Turing devised an incredibly simple kind of computing mechanism (what we now call, in his honor, a *Turing machine*). This is an idealized machine, not a real one. What makes a Turing machine idealized is that it consists of an infinitely long piece of tape divided into cells. The point of the tape being infinitely long is so that the machine will not have any storage limitations. A Turing machine is like a computer with an infinitely large hard disk. Turing did not think that a Turing machine would ever have to deal with infinitely long strings of symbols. He just wanted it to be able to deal with arbitrarily long, but still finite, strings of symbols.

Each of the cells of the Turing tape can be either blank or contain a single symbol. The Turing machine contains a machine head. The tape runs through the machine head, with a single cell under the head at a given moment. This allows the head to read the symbol the cell contains. The machine head can also carry out a limited number of operations on the cell that it is currently scanning. It can:

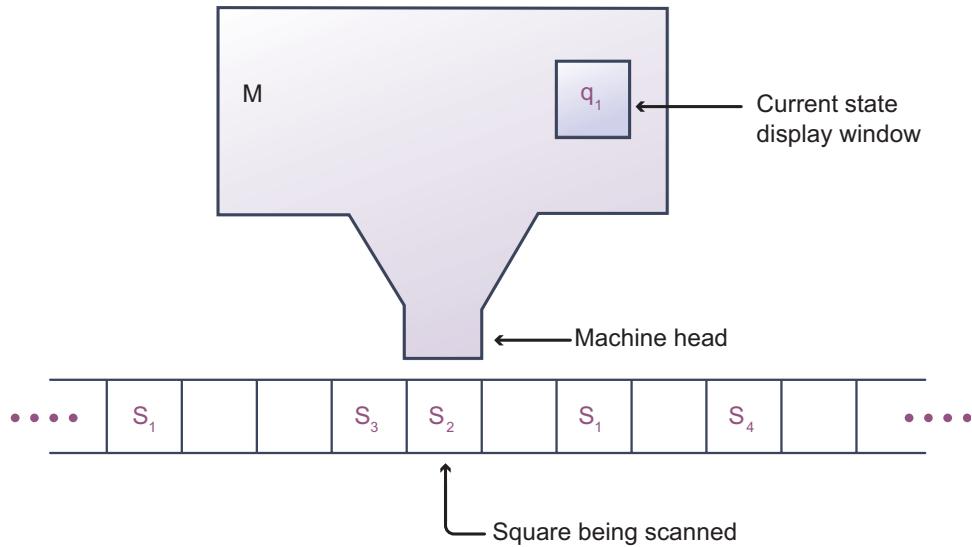
- delete the symbol in the cell
- write a new symbol in the cell
- move the tape one cell to the left
- move the tape one cell to the right

Any individual Turing machine has a set of instructions (its *machine table*). The machine can be in any one of a (finite) number of different states. The machine table determines what the Turing machine will do when it encounters a particular symbol in a particular cell, depending upon which internal state it is in. Figure 1.4 is a schematic representation of a Turing machine.

The beauty of a Turing machine is that its behavior is entirely determined by the machine table, its current state, and the symbol in the cell it is currently scanning. There is no ambiguity and no room for the machine to exercise “intuition” or “judgment.” It is, in fact, purely mechanical in exactly the way required for an algorithm.

Turing did not actually build a Turing machine. (It is difficult to build a machine with an infinitely long piece of tape!) But he showed how Turing machines could be specified mathematically. The machine table of a Turing machine can be represented as a sequence of numbers. This allowed him to prove mathematical results about Turing machines. In particular, it allowed him to prove that there is a special kind of Turing machine, a *universal Turing machine*, that can run any specialized Turing machine. The universal Turing machine can take as input a program specifying any given specialized Turing program. It is the theoretical precursor (with unlimited storage) of the modern-day general-purpose digital computer.

Turing's paper contained a subtle proof that the Halting Problem cannot be solved. It was also significant for articulating what we now call the *Church–Turing thesis* (in recognition of the contribution made by the logician Alonzo Church). According to the Church–Turing thesis, anything that can be done in mathematics by an algorithm can



**Figure 1.4** Schematic representation of a Turing machine. (Adapted from Cutland 1980)

be done by a Turing machine. Turing machines are computers that can compute anything that can be algorithmically computed.

What Turing contributed to the early development of cognitive science (although at the time his work was little known and even less appreciated) was a model of computation that looked as if it might be a clue to how information could be processed by the mind. As theorists moved closer to the idea that cognition involves processing information it was an easy step to think about information processing as an algorithmic process along the lines analyzed by Turing – a step that became even easier in the light of the huge advances that were made in designing and building digital computers (which, if the Church–Turing thesis is true, are essentially large and fast Turing machines) during and after the Second World War.



**Exercise 1.5** Explain in your own words why the Church–Turing thesis entails that any computer running a program is simply a large and fast Turing machine.

## 1.3 Linguistics and the Formal Analysis of Language

The study of language played a fundamental role in the prehistory of cognitive science. On the one hand, language use is a paradigm of the sort of hierarchically organized complex behavior that Lashley was talking about. On the other hand, the emergence of transformational linguistics and the formal analysis of *syntax* (those aspects of language use that have to do with how words can be legitimately put together to form sentences) provided a very clear example of how to analyze, *in algorithmic terms*, the bodies of information that might underlie certain very basic cognitive abilities (such as the ability to speak and understand a language).

In retrospect we can identify one crucial landmark as the publication in 1957 of *Syntactic Structures* by Noam Chomsky, unquestionably the father of modern linguistics and a hugely important figure in the development of cognitive science. The transformational grammar proposed by Chomsky (and subsequently much modified by Chomsky and others) reflects some of the basic ideas covered earlier in this chapter.

## The Structure of Language: Chomsky's *Syntactic Structures* (1957)

Chomsky's book is widely held to be the first example of a linguist proposing an explanatory theory of *why* languages work the way they do (as opposed to simply describing and classifying *how* they work). Chomsky was interested not in mapping the differences between different languages and in describing their structure, but rather in providing a theoretical account of why they have the structure that they do. Crucial to his approach is the distinction between the *deep structure* of a sentence (as given by what Chomsky calls a *phrase structure grammar*) and its *surface structure* (the actual organization of words in a sentence, derived from the deep structure according to the principles of transformational grammar).

The deep structure, or phrase structure, of a sentence is simply how it is built up from basic constituents (syntactic categories) according to basic rules (phrase structure rules). We only need a small number of basic categories to specify the phrase structure of a sentence. These are the familiar parts of speech that we all learn about in high school – nouns, verbs, adjectives, and so on. Any grammatical sentence (including those that nobody is ever likely to utter) is made up of these basic parts of speech combined according to basic phrase structure rules (such as the rule that every sentence is composed of a verb phrase and a noun phrase).

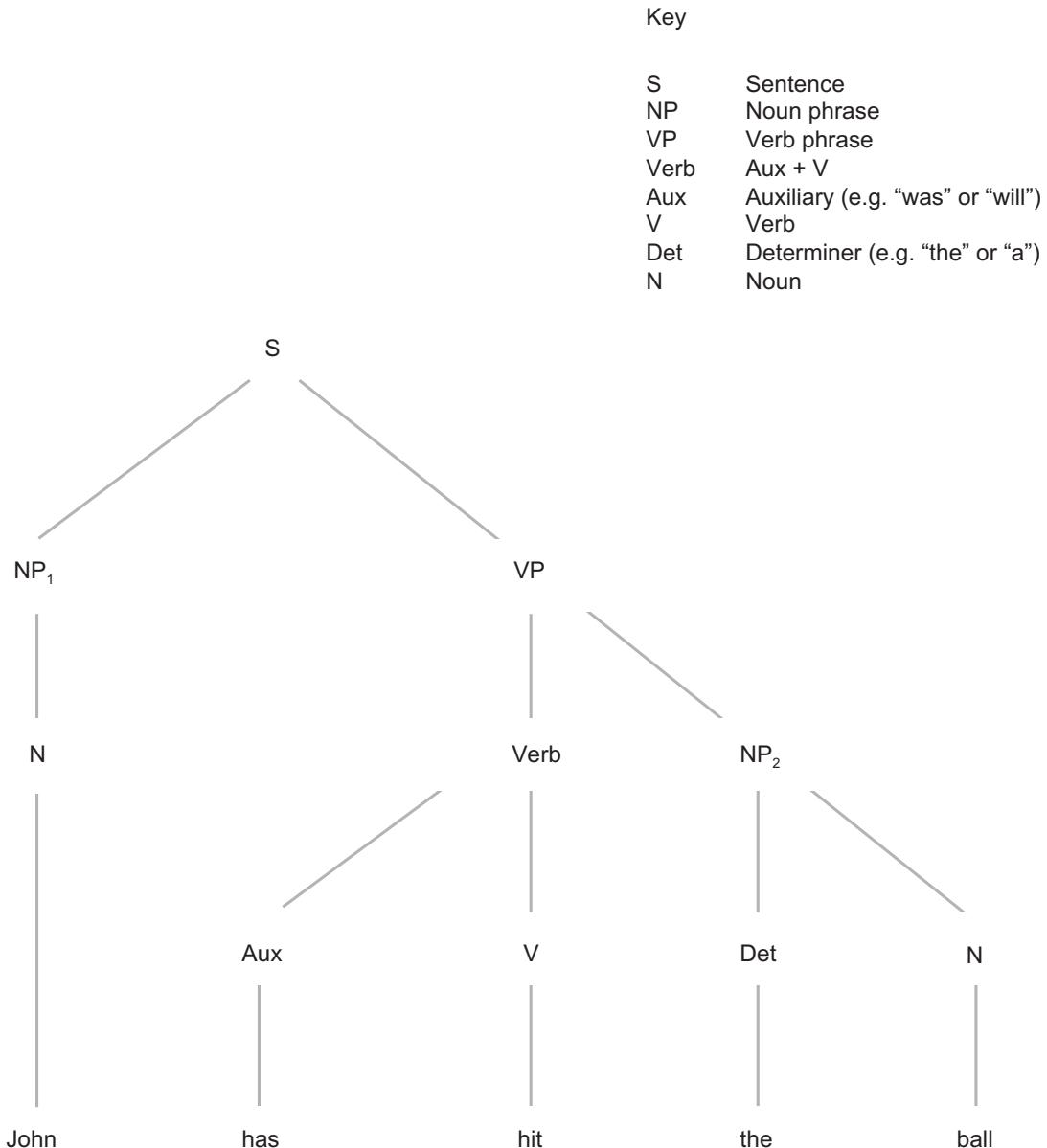
In Figure 1.5 we see how these basic categories can be used to give a phrase structure tree of the sentence “John has hit the ball.” The phrase structure tree is easy to read, with a bit of practice. Basically, you start at the top with the most general characterization. As you work your way down the tree the structure of the sentence becomes more finely articulated, so that we see which words or combinations of words are doing which job.

Analyzing sentences in terms of their phrase structure is a powerful explanatory tool. There are pairs of sentences that have very different phrase structures, but are clearly very similar in meaning. Think of “John has hit the ball” and “The ball has been hit by John.” In most contexts these sentences are equivalent and interchangeable, despite having very different phrase structures. Conversely, there are sentences with superficially similar phrase structures that are plainly unrelated. Think of “Susan is easy to please” and “Susan is eager to please.”



**Exercise 1.6** Explain in your own words the difference between these two sentences. Why are their phrase structures different?

The basic aim of transformational grammar is to explain the connection between sentences of the first type and to explain the differences between sentences of the second



**Figure 1.5** A sample phrase structure tree for the sentence "John has hit the ball." The abbreviations in the diagram are explained in the key.

type. This is done by giving principles that state the acceptable ways of transforming deep structures. This allows linguists to identify the transformational structure of a sentence in terms of its transformational history.

The transformational principles of transformational grammar are examples of *algorithms*. They specify a set of procedures that operate upon a string of symbols to convert it into a different string of symbols. So, for example, our simple phrase structure grammar might be extended to include an active–passive transformation rule that takes the following form (look at the key in Figure 1.5 for the translation of the symbols):

$$\begin{aligned} \text{NP}_1 + \text{Aux} + \text{V} + \text{NP}_2 \\ \Rightarrow \\ \text{NP}_2 + \text{Aux} + \text{been} + \text{V} + \text{by} + \text{NP}_1 \end{aligned}$$

This transforms the string “John + has + hit + the + ball” into the string “the + ball + has + been + hit + by + John.” And it does so in a purely mechanical and algorithmic way.



**Exercise 1.7** Write out an algorithm that carries out the active–passive transformation rule. Make sure that your algorithm instructs the person/machine following it what to do at each step.

What’s more, when we look at the structure of the passive sentence “The ball has been hit by John” we can see it as illustrating precisely the sort of hierarchical structure to which Lashley drew our attention. This is a characteristic of languages in general. They are hierarchically organized.

So, in thinking about how they work, transformational grammar brings together two very fundamental ideas. The first idea is that a sophisticated, hierarchically organized, cognitive ability, such as speaking and understanding a language, involves stored bodies of information (information about phrase structures and transformation rules). The second idea is that these bodies of information can be manipulated algorithmically.

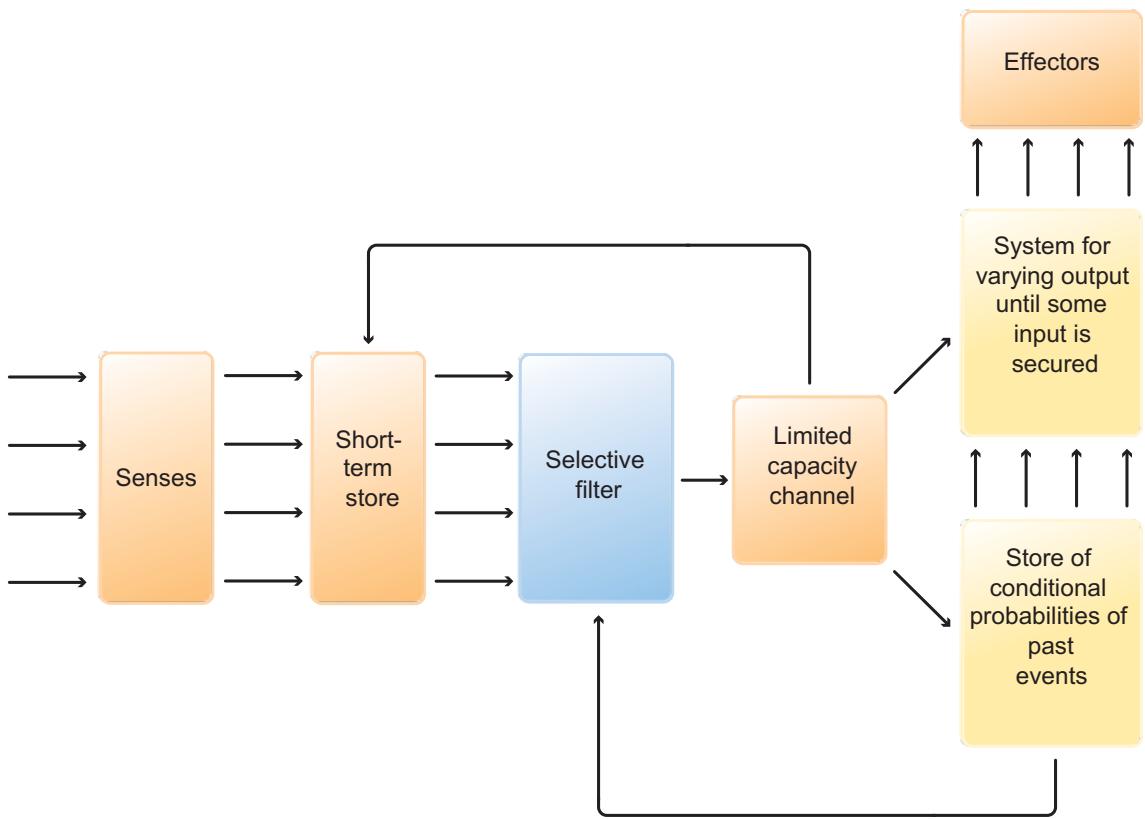
## 1.4

## Information-Processing Models in Psychology

In the late 1950s the idea that the mind works by processing information began to take hold within psychology. This new development reflected a number of different influences. One of these was the emergence of information theory in applied mathematics. Rather unusually in the history of science, the emergence of information theory can be pinned down to a single event – the publication of an article entitled “A mathematical theory of communication” by Claude E. Shannon in 1948. Shannon’s paper showed how information can be measured, and he provided precise mathematical tools for studying the transmission of information.

These tools (including the idea of a *bit* as a measure of information) proved very influential in psychology, and for cognitive science more generally. We can illustrate how information-processing models became established in psychology through two very famous publications from the 1950s.

The first, George Miller’s article “The magical number seven, plus or minus two: Some limits on our capacity for processing information,” used the basic concepts of information theory to identify crucial features of how the mind works. The second, Donald Broadbent’s 1954 paper “The role of auditory localization in attention and memory span,” presented two influential experiments that were crucial in Broadbent’s later putting forward, in his 1958 book *Perception and Communication*, one of the first information-processing models in psychology. The type of flowchart model that Broadbent proposed (as illustrated in Figure 1.6) has become a standard way for cognitive scientists to describe and explain different aspects of cognition.



**Figure 1.6** Donald Broadbent's (1958) model of selective attention.

## How Much Information Can We Handle? George Miller's "The Magical Number Seven, Plus or Minus Two" (1956)

The tools of information theory can be applied to the study of the mind. One of the basic concepts of information theory is the concept of an information channel. In abstract terms, an information channel is a *medium* that transmits information from a *sender* to a *receiver*. A telephone cable is an information channel. So is the radio frequency on which a television station broadcasts. Perceptual systems are themselves information channels. Vision, for example, is a medium through which information is transmitted from the environment to the perceiver. So are audition (hearing) and olfaction (smell). Thinking about perceptual systems in this way gave Miller and other psychologists a new set of tools for thinking about experiments on human perception.

Miller's article drew attention to a wide range of evidence suggesting that human subjects are really rather limited in the *absolute judgments* that they can make. An example of an absolute judgment is naming a color, or identifying the pitch of a particular tone – as opposed to relative judgments, such as identifying which of two colors is the darker, or which of two tones is higher in pitch.

In one experiment reported by Miller, subjects are asked to assign numbers to the pitches of particular tones and then presented with sequences of tones and asked to

identify them in terms of the assigned numbers. So, for example, if you assigned “1” to middle C, “2” to the first E above middle C, and “3” to the first F# and then heard the sequence E-C-C-F#-E, the correct response would be 2-1-1-3-2.

When the sequence is only one or two tones long, subjects never make mistakes. But performance falls off drastically when the sequence is six or more tones long. A similar phenomenon occurs when we switch from audition to vision and ask subjects to judge the size of squares or the length of a line. Here too there seems to be an upper bound on the number of distinct items that can be processed simultaneously.

Putting these (and many other) experimental results into the context of information theory led Miller to propose that our sensory systems are all information channels with roughly the same *channel capacity* (where the channel capacity of an information channel is given by the amount of information it can reliably transmit). In these cases the perceiver’s capacity to make absolute judgments is an index of the channel capacity of the information channel that she is using.

What Miller essentially did was propose an *information-processing bottleneck*. The human perceptual systems, he suggested, are information channels with built-in limits. These information channels can only process around seven items at the same time (or, to put it in the language of information theory, their channel capacity is around 3 bits; since each bit allows the system to discriminate 2 pieces of information,  $n$  bits of information allow the system to discriminate  $2^n$  pieces of information and 7 is just under  $2^3$ ).

At the same time as identifying these limits, Miller identified ways of working round them. One way of increasing the channel capacity is to *chunk* information. We can relabel sequences of numbers with single numbers. A good example (discussed by Miller) comes when we use decimal notation to relabel numbers in binary notation. We can pick out the same number in two different ways – with the binary expression 1100100, for example, or with the decimal expression 100. If we use binary notation then we are at the limits of our visual channel capacity. If we use decimal notation then we are well within those limits. As Miller pointed out, to return to a theme that has emerged several times already, natural language is the ultimate chunking tool.



**Exercise 1.8** Think of an informal experiment that you can do to illustrate the significance of chunking information.

## The Flow of Information: Donald Broadbent's "The Role of Auditory Localization in Attention and Memory Span" (1954) and *Perception and Communication* (1958)

Miller’s work drew attention to some very general features of how information is processed in the mind, but it had little to say about the details of how that information processing takes place. The experiments reported and analyzed by Miller made plausible the idea that the senses are information channels with limited capacity. The obvious next step was to think about how those information channels actually work. One of the first models of how sensory information is processed was developed by the British



psychologist Donald Broadbent in his 1958 book *Perception and Communication*. As with Miller, the impetus came from experiments in the branch of psychology known as *psychophysics*. This is the branch of psychology that studies how subjects perceive and discriminate physical stimuli.

We can appreciate what is going on by thinking about the so-called *cocktail party phenomenon*. When at a cocktail party, or any other social gathering, we can often hear many ongoing and unrelated conversations. Somehow we manage to focus only on the one we want to listen to. How do we manage this? How do we screen out all the unwanted sentences that we hear? It is plain that we only *attend* to some of what we hear. Auditory attention is selective. There is nothing peculiar to audition here, of course. The phenomenon of *selective attention* occurs in every sense modality.

Broadbent studied auditory attention by using *dichotic listening experiments*, in which subjects are presented with different information in each ear. The experiments reported in his paper "The role of auditory localization in attention and memory span" involved presenting subjects with a string of three different stimuli (letters or digits) in one ear, while simultaneously presenting them with a different string in the other ear. The subjects were asked to report the stimuli in any order. Broadbent found that they performed best when they reported the stimuli ear by ear – that is, by reporting all three presented to the left ear first, followed by the three presented to the right ear. This, and other findings, were explained by the model that he subsequently developed.

The basic features of the model are illustrated in Figure 1.6. Information comes through the senses and passes through a short-term store before passing through a selective filter. The selective filter screens out a large portion of the incoming information, selecting some of it for further processing. This is what allows us selectively to attend to only a portion of what is going on around us in the cocktail party. Only information that makes it through the selective filter is semantically interpreted, for example. Although people at cocktail parties can hear many different conversations at the same time, many experiments have shown that they have little idea of what is said in the conversations that they are not attending to. They hear the words, but do not extract their meaning.

Broadbent interpreted the dichotic listening experiments as showing that we can only attend to a single information channel at a time (assuming that each ear is a separate information channel) – and that the selection between information channels is based purely on physical characteristics of the signal. The selection might be based on the physical location of the sound (whether it comes from the left ear or the right ear, for example), or on whether it is a man's voice or a woman's voice.

The selective filter does not work by magic. As the diagram shows, the selective filter is "programmed" by another system that stores information about the relative likelihoods of different events. We are assuming that the system is pursuing a goal. What is programming the selective filter is information about the sorts of things that have led to that goal being satisfied in the past. Information that makes it through the selective filter goes into what Broadbent calls the limited capacity channel. Information that is filtered out is assumed to decay quickly. From the limited capacity channel information can go either into the long-term store, or on to further processing and eventually into action, or it can be recycled back into the short-term store (to preserve it if it is in danger of being lost).

We can see how Broadbent's model characterizes what is going on in the cocktail party phenomenon. The stream of different conversations arrives at the selective filter. If my goal, let us say, is to strike up a conversation with Dr. X (who is female), then the selective filter might be attuned in the first instance to female voices. The sounds that make it through the selective filter are the sounds of which I am consciously aware. They can provide information that can be stored and perhaps eventually feed back into the selective filter. Suppose that I "tune into" a conversation that I think involves Dr. X but where the female voice turns out to belong to Mrs. Z, then the selective filter can be instructed to filter out Mrs. Z's voice.



**Exercise 1.9** Give an example in your own words of selective attention in action. Incorporate as many different aspects of Broadbent's model as possible.



## 1.5

## Connections and Points of Contact

This chapter has surveyed some crucial episodes in the prehistory of cognitive science. You should by now have a sense of exciting innovations and discoveries taking place in very different areas of intellectual life – from experiments on rats in mazes to some of the most abstract areas of mathematics, and from thinking about how we navigate cocktail parties to analyzing the deep structure of natural language. As we have looked at some of the key publications in these very different areas, a number of fundamental ideas have kept recurring.

The most basic concept that has run through the chapter is the concept of information. Tolman's latent-learning experiments seemed to many to show that animals (including of course human animals) are capable of picking up information without any reinforcement taking place. The rats wandering unrewarded through the maze were picking up and storing information about how it was laid out – information that they could subsequently retrieve and put to work when there was food at stake.

Chomsky's approach to linguistics exploits the concept of information in a very different way. His *Syntactic Structures* pointed linguists toward the idea that speaking and understanding natural languages depends upon information about sentence structure – about the basic rules that govern the surface structure of sentences and about the basic transformation principles that underlie the deep structure of sentences.

In the work of the psychologists Miller and Broadbent we find the concept of information appearing in yet another form. Here the idea is that we can understand perceptual systems as information channels and use the concepts of information theory to explore their basic structure and limits.

Hand in hand with the concept of information goes the concept of representation. Information is everywhere, but in order to use it organisms need to represent it. Representations will turn out to be the basic currency of cognitive science, and we have seen a range of very different examples of how information is represented in this chapter.

Tolman's place-learning experiments introduced the idea that organisms have cognitive maps representing the spatial layout of the environment. These maps are representations of the environment. Turing machines incorporate a very different type of representation. They



represent the instructions for implementing particular algorithms in their machine table. In a similar vein, Chomsky suggested that important elements of linguistic understanding are represented as phrase structure rules and transformational rules. And Miller showed how representing information in different ways (in terms of different types of chunking, for example) can affect how much information we are able to store in memory.

Information is not a static commodity. Organisms pick up information. They adapt it, modify it, and use it. In short, organisms engage in *information processing*. The basic idea of information processing raises a number of questions. One might wonder, for example, about the *content* of the information that is being processed. What an organism does with information depends upon how that information is encoded.

We saw some of the ramifications of this in Tolman's place-learning experiments. The difference between place learning and response learning is a difference in how information about location is encoded. In response learning, information about location is encoded in terms of the movements that an organism might make to reach that location. In place learning, in contrast, information about location is encoded in terms of the location's relation to other locations in the environment.

Even once we know how information is encoded, there remain questions about the mechanics of information processing. How does it actually work? We can see the germ of a possible answer in Turing's model of computation. Turing machines illustrate the idea of a purely mechanical way of solving problems and processing information. In one sense Turing machines are completely unintelligent. They blindly follow very simple instructions. And yet, if the Church-Turing thesis is warranted, they can compute anything that can be algorithmically computed. And so, in another sense, it would be difficult to be more intelligent than a Turing machine.

If the basic assumptions of transformational linguistics are correct, then we can see one sphere in which the notion of an algorithm can be applied. The basic principles that transform sentences (that take a sentence from its active to its passive form, for example, or that transform a statement into a question) can be thought of as mechanical procedures that can in principle be carried out by a suitably programmed Turing machine (once we have found a way of numerically coding the basic categories of transformational grammar).

A final theme that has emerged from the authors we have studied is the idea that information processing is done by dedicated and specialized systems. This idea comes across most clearly in Broadbent's model of selective attention. Here we see a complex information-processing task (the task of making sense of the vast amounts of information picked up by the hearing system) broken down into a number of simpler tasks (such as the task of selecting a single information channel, or the task of working out what sentences mean). Each of these information-processing tasks is performed by dedicated systems, such as the selective filter or the semantic processing system.

One powerful idea that emerges from Broadbent's model of selective attention is the idea that we *can understand how a cognitive system as a whole works by understanding how information flows through the system*. What Broadbent offered was a flowchart showing the different stages that information goes through as it is processed by the system. Many psychologists and cognitive scientists subsequently took this type of information-processing flowchart to be a paradigm of how to explain cognitive abilities.

In the next chapter we will look at how some of these ideas were put together in some of the classic theories and models of early cognitive science.



## Summary

This chapter has surveyed five of the most important precursors of what subsequently became known as cognitive science. Cognitive science emerged when experimentalists and theoreticians began to see connections between developments in disciplines as diverse as experimental psychology, theoretical linguistics, and mathematical logic. These connections converge on the idea that cognition is a form of information processing and hence that we can understand how the mind works and how organisms negotiate the world around them by understanding how information about the environment is represented, transformed, and exploited.

## Checklist

### Important Developments Leading Up to the Emergence of Cognitive Science

- (1) The reaction against behaviorism in psychology
- (2) Theoretical models of computation from mathematical logic
- (3) Systematic analysis of the structure of natural language in linguistics
- (4) The development of information-processing models in psychology

### Central Themes of the Chapter

- (1) Even very basic types of behavior (such as the behavior of rats in mazes) seems to involve storing and processing information about the environment.
- (2) Information relevant to cognition can take many forms – from information about the environment to information about how sentences can be constructed and transformed.
- (3) Perceptual systems can be viewed as information channels and we can study both: (a) the very general properties of those channels (e.g., their channel capacity) (b) the way in which information flows through those channels.
- (4) Mathematical logic and the theory of computation shows us how information processing can be mechanical and algorithmic.
- (5) Much of the information-processing that goes on in the mind takes place below the threshold of awareness.

## Further Reading

The story of how cognitive science emerged is told in Gardner's *The Mind's New Science* (1985). Flanagan's *The Science of the Mind* (1991) goes further back into the prehistory of cognitive science and psychology, as do the papers in Brook 2007. Margaret Boden's two-volume *Mind as Machine: A History of Cognitive Science* (2006) is detailed, placing most emphasis on computer science and artificial intelligence. Abrahamsen and Bechtel's chapter in Frankish and Ramsey 2012 provides a concise summary of the history of cognitive science. Going back still further, histories of psychology, such as Hergenhahn and Helnley 2013, typically start with the ancient Greeks.



The basic principles of classical and operant conditioning are covered in standard textbooks to psychology, such as Gazzaniga, Halpern, and Heatherton 2011, Plotnik and Kouyoumdjian 2010, and Kalat 2010. For a survey of contemporary research on both types of conditioning, see McSweeney and Murphy 2014. Watson's article "Psychology as the behaviorist views it" is a classic behaviorist manifesto (Watson 1913). It can be found in the online resources. Tolman's article "Cognitive maps in rats and men" (1948) gives an accessible introduction to many of his experiments and is also in the online resources. Gallistel 1990 is a very detailed and sophisticated presentation of a computational approach to animal learning.

Turing's paper on undecidable propositions (Turing 1936) will defeat all but graduate students in mathematical logic (but see Petzold 2008 for a book-length explanation aimed at the general reader). His paper "Computing machinery and intelligence" (Turing 1950) is a much more accessible introduction to his thoughts about computers. There are several versions online, the best of which are included in the online resources. Hodges 2014 is a new edition of the classic biography of Turing, which inspired the film *The Imitation Game*. Martin Davis has written two popular books on the early history of computers, *Engines of Logic: Mathematicians and the Origin of the Computer* (2001) and *The Universal Computer: The Road from Leibniz to Turing* (2000). Copeland 1993 gives a more technical, but still accessible, account of Turing machines and the Church–Turing thesis. Millican and Clark 1996 is a collection of papers on Turing's legacy. A more general article illustrating the algorithmic nature of information processing is Schyns, Gosselin, and Smith 2008.

At more or less the same time as Turing was working on the mathematical theory of computation, the neurophysiologist Warren McCulloch and logician Walter Pitts were collaborating on applying rather similar ideas about computation directly to the brain. Their paper "A logical calculus of the ideas immanent in nervous activity" (McCulloch and Pitts 1943) was influential at the time, particularly in the early development of digital computers, but is rarely read now. It is reprinted in Cummins and Cummins 2000. An accessible survey of their basic ideas can be found in Anderson 2003. See also chapter 2 of Arbib 1987 and Piccinini 2004, as well as Schlatter and Aizawa 2008.

Most people find Chomsky's *Syntactic Structures* pretty hard going. Linguistics tends to be technical, but Chomsky's article "Linguistics and philosophy," reprinted in Cummins and Cummins 2000, contains a fairly informal introduction to the basic distinction between surface structure and deep structure. Chapter 2 of Newmeyer 1986 is a good and accessible introduction to the Chomskyan revolution. More details can be found in standard textbooks, such as Cook and Newsom 2007, Isac and Reiss 2013, and O'Grady et al. 2010. Chomsky's rather harsh review of B. F. Skinner's book *Verbal Behavior* (Chomsky 1959) is often described as instrumental in the demise of radical behaviorism – and hence in bringing about the so-called cognitive revolution. The review is reprinted in many places and can be found in the online resources. Pinker 1994 presents a broadly Chomskyan perspective on language for a general audience.

The cocktail party phenomenon was first introduced in Cherry 1953. A concise summary of the cocktail party phenomenon can be found in McDermott 2009. Miller's 1956 article is widely available and is included in the online resources. Broadbent's model of selective attention was the first in a long line of models. These are reviewed in standard textbooks. See, for example, chapter 5 of Gleitman, Fridlund, and Reisberg 2010. Christopher Mole's chapter on attention in Margolis, Samuels, and Stich 2012 summarizes Broadbent's influence as well as recent departures from Broadbent. Driver 2001 is an article-length survey of theories of selective attention in the twentieth century.





## CHAPTER TWO

# The Discipline Matures: Three Milestones

### OVERVIEW 37

- 2.1 Language and Micro-worlds** 38  
Natural Language Processing: Winograd,  
*Understanding Natural  
Language* (1972) 39  
SHRDLU in Action 41
- 2.2 How Do Mental Images  
Represent?** 47

Mental Rotation: Shepard and Metzler,  
“Mental Rotation of Three-Dimensional  
Objects” (1971) 48

Information Processing in Mental Imagery 50

- 2.3 An Interdisciplinary Model of Vision** 53  
Levels of Explanation: Marr’s  
*Vision* (1982) 53  
Applying Top-Down Analysis to the Visual  
System 55



## Overview

Chapter 1 explored the prehistory of cognitive science in the first half of the twentieth century. In this second chapter of this selective historical survey we will look closely at three milestones in the development of cognitive science. In each of them we start to see some of the theoretical ideas canvassed in the previous section being combined and applied to understanding specific cognitive systems and cognitive abilities.

Section 2.1 looks at a powerful and influential computer model of what it is to understand a natural language. Terry Winograd’s computer model SHRDLU illustrates how grammatical rules might be represented in a cognitive system and integrated with other types of information about the environment. SHRDLU’s programming is built around specific procedures that carry out fairly specialized information-processing tasks in an algorithmic (or at least quasi-algorithmic way).

The idea that the digital computer is the most promising model for understanding the mind was at the forefront of cognitive science in the 1960s and 1970s. But even in the 1970s it was under pressure. Section 2.2 looks at the debate on the nature of mental imagery provoked by some very influential experiments in cognitive psychology. These experiments seemed to many theorists to

show that some types of cognitive information processing involve forms of representation very different from how information is represented in, and manipulated by, a digital computer.

The third section introduces what many cognitive scientists still consider to be cognitive science's greatest single achievement – the theory of early visual processing developed by David Marr. Marr's theory of vision was highly interdisciplinary, drawing on mathematics, cognitive psychology, neuroscience, and the clinical study of brain-damaged patients and it was built on a hierarchy of different levels for studying cognition that was for a long time taken to define the method of cognitive science.



## 2.1

## Language and Micro-worlds

The human ability to speak and understand natural language is one of our most sophisticated cognitive achievements. We share many types of cognitive ability with nonlinguistic animals. Many cognitive scientists assume, for example, that there are significant continuities between human perceptual systems and those of the higher primates (such as chimpanzees and macaque monkeys), which is why much of what we know about the neural structure of the human perceptual system is actually derived from experiments on monkeys. (More on this in Chapters 3 and 9.) And there is powerful evidence that prelinguistic infants are capable of representing and reasoning about their physical and social environment in comparatively sophisticated ways. (See Chapter 11 for more details.)

Nonetheless, just as in human development (*ontogeny*) there is a cognitive explosion that runs more or less in parallel with the acquisition of language, much of what distinguishes humans from other animals is intimately bound up with our linguistic abilities. Language is far more than a tool for communication. It is a tool for thinking. Without language there would be no science and no mathematics. Language allows us to engage in incredibly sophisticated types of coordinated behavior. It underpins our political and social structures. In many ways, *Homo linguisticus* would be a more accurate name than *Homo sapiens*.

Unsurprisingly, then, the study of natural language has always been at the center of cognitive science. If cognitive scientists want to understand the human mind then they have to confront the fundamental challenge posed by our understanding of natural language. As we saw in the last chapter, Chomsky's diagnosis of what he saw as the insuperable challenges facing a behaviorist account of language was very important in setting the stage for the cognitive revolution. So too was the discovery, also due to Chomsky, of ways of describing the underlying structures that lie beneath the patterns of surface grammar.

But Chomsky's transformational linguistics has relatively little to say about how linguistic understanding actually takes place. It is one thing to describe the abstract structure of human language and quite another to explain how human beings can master that abstract structure. What Chomsky's work tells us (if it is indeed the correct way to think about the deep structure of language) is what we know when we understand a language. It tells us what we have to know. But it has nothing to say about how that knowledge is stored or how it is used.



## Natural Language Processing: Winograd, Understanding Natural Language (1972)

The first study that we examine in this chapter confronts this challenge head on. One way of trying to model how we store and use linguistic knowledge is to build a machine that is capable of some form of linguistic understanding. The early days of *artificial intelligence* (AI) saw a number of attempts to write computer programs that could engage in some very elementary forms of conversational exchanges, but none of these programs was capable of anything that really resembled linguistic understanding.

The aim of programs such as ELIZA (written by Joseph Weizenbaum in 1966) was to simulate human conversation. The basic idea behind ELIZA (which, depending upon who one asks, was either based upon or intended to parody typical conversational exchanges between psychotherapists and their patients) was to create the illusion of conversation by rephrasing statements as questions and by programming the computer to give certain fixed responses where this is not possible. A sample “conversation” is given in Box 2.1.

Although ELIZA is said to have fooled a number of people into thinking that it was a human (including the unknowing participant in the conversation recorded in the box) nobody has ever suggested that it displays anything like a genuine understanding of

### BOX 2.1 A Conversation with ELIZA

ELIZA's responses are in uppercase.

Men are all alike  
IN WHAT WAY?

They are always bugging us about something or other  
CAN YOU THINK OF A SPECIFIC EXAMPLE?

Well, my boyfriend made me come here  
YOUR BOYFRIEND MADE YOU COME HERE?

He says I'm depressed much of the time  
I'M SORRY TO HEAR THAT YOU ARE DEPRESSED

It's true. I am unhappy  
DO YOU THINK COMING HERE WILL HELP YOU NOT TO BE UNHAPPY?

I need some help. That much seems certain  
WHAT WOULD IT MEAN TO YOU IF YOU GOT SOME HELP?

Perhaps I could learn to get along with my mother  
TELL ME MORE ABOUT YOUR FAMILY

(from Boden 1977: 106–7)

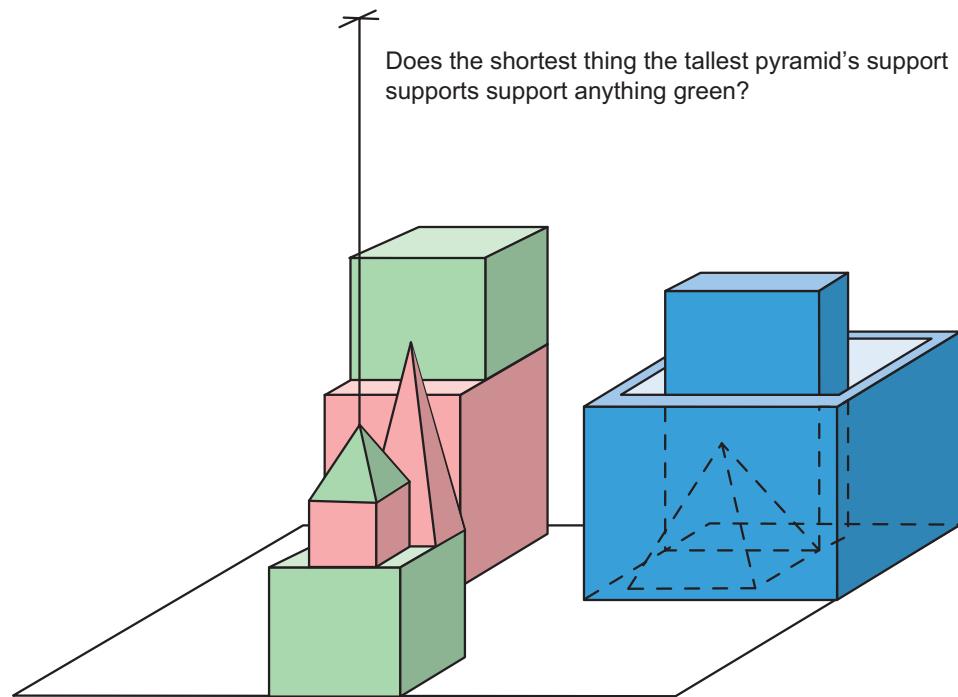
language. ELIZA does not in any sense analyze the syntactic structure or the meaning of the sentences that it encounters. It is simply programmed to respond to certain cues by making one of a small set of responses. Nor could ELIZA use the conversations in which it engaged to report on or navigate its environment. *Chatbot* programs such as ELIZA are interesting for many reasons, but not as serious models of how we understand and use language.



**Exercise 2.1** Explain in your own words what you think we can learn from programs such as ELIZA. Is it important that a person might be fooled by ELIZA into thinking that we were communicating with another human being?

Terry Winograd's program SHRDLU, initially presented in his 1970 doctoral dissertation at MIT, was one of the first attempts to write a program that was not just trying to simulate conversation, but that was capable of using language to report on its environment, to plan actions, and to reason about the implications of what is being said to it.

One of the distinctive features of SHRDLU is that it is programmed to deal with a very limited *micro-world* (as opposed to being a general-purpose language program, which is what ELIZA and other chatterbot programs are, in their very limited ways). The SHRDLU micro-world is very simple. It consists simply of a number of colored blocks, colored pyramids, and a box, all located on a tabletop, as illustrated in Figure 2.1. (The micro-world is a virtual micro-world, it should be emphasized. Everything takes place on a computer screen.)



**Figure 2.1** A question for SHRDLU about its virtual micro-world. (Adapted from Winograd 1972)



SHRDLU is capable of various actions in the micro-world, which it can carry out through a (virtual) robot arm. It can pick up the blocks and pyramids, move them around, and put them in the box. Corresponding to the simplicity of the micro-world, SHRDLU's language is relatively simple. It only has the tools to talk about what is going on in the micro-world.

SHRDLU was very important in the development of cognitive science, for three main reasons. First, it gave a powerful illustration of how abstract rules and principles such as those in the sort of grammar that we might find in theoretical linguistics could be practically implemented. If we assume that a speaker's understanding of language is best understood as a body of knowledge, then SHRDLU provided a model of how that knowledge could be *represented* by a cognitive system and how it could be *integrated* with other, more general, forms of knowledge about the environment.

Second, SHRDLU illustrates the general approach of trying to understand and model cognitive systems by breaking them down into distinct components, each carrying out a specific information-processing task. One of the many interesting things about SHRDLU is that these distinct components are not completely self-contained. The separate processing systems collaborate in solving information-processing problems. There is *cross-talk* between them, because the programs for each processing system allow it to consult other processing systems at particular moments in the computation.

Finally, the SHRDLU is based on the fundamental assumption that understanding language is an *algorithmic* process. In Winograd's own words, "All language use can be thought of as a way of activating procedures within the hearer" (1973: 104). Each component system is essentially made up of a vast number of procedures that work algorithmically to solve very specific problems. The system as a whole works because of how these procedures are linked up and embedded within each other.

## SHRDLU in Action

As is often the case in so-called *classical cognitive science*, the best way to understand what is going on in SHRDLU is to work from the top down – to start by looking at the general overall structure and then drill down into the details. Strictly speaking, SHRDLU consists of twelve different systems. Winograd himself divides these into three groups. Each group carries out a specific job. The particular jobs that Winograd identifies are not particularly surprising. They are exactly the jobs that one would expect any language-processing system to carry out.

- 1 *The job of syntactic analysis:* SHRDLU needs to be able to "decode" the grammatical structure of the sentences that it encounters. It needs to be able to identify which units in the sentence are performing which linguistic function. In order to *parse* any sentence, a language user needs to work out which linguistic units are functioning as nouns (i.e., are picking out objects) and which are functioning as verbs (i.e., characterizing events and processes).
- 2 *The job of semantic analysis:* Understanding a sentence involves much more than decoding its syntactic structure. The system also needs to assign meanings to the individual words in a way that reveals what the sentence is stating (if it is a statement), or requesting (if it is a request). This takes us from *syntax* to *semantics*.

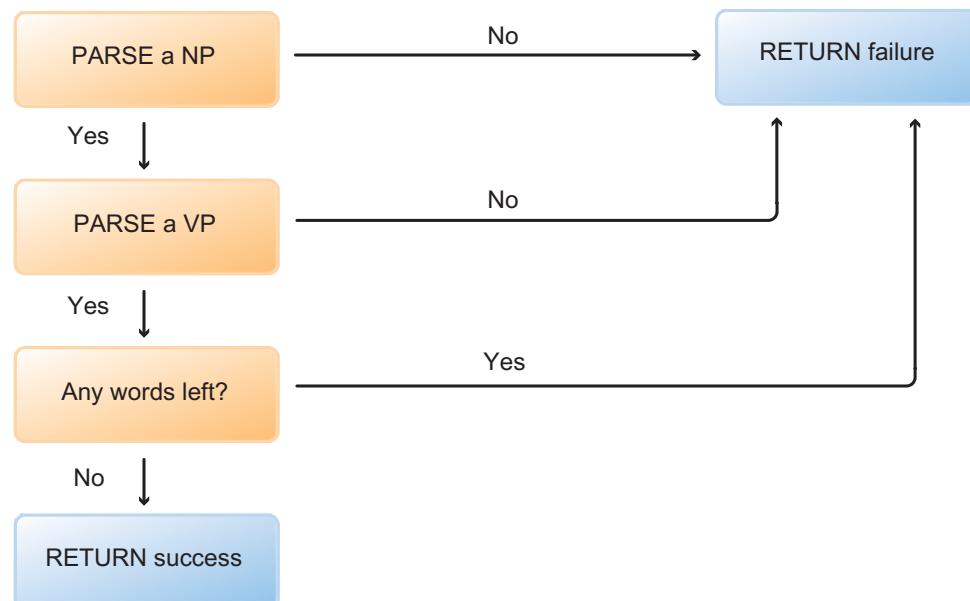
- 3** *The job of integrating the information acquired with the information the system already possesses:* The system has to be able to explore the implications of what it has just learned for the information it already has. Or to call upon information it already has in order to obey some command, fulfill a request, or answer a question. These all require ways of deducing and comparing the logical consequences of stored and newly acquired information.

We can identify distinct components for each of these jobs – the *syntactic system*, the *semantic system*, and the *cognitive-deductive system*. Winograd does not see these as operating in strict sequence. It is not the case that the syntactic system does its job producing a syntactic analysis, and then hands that syntactic analysis over to the semantic system, which plugs meanings into the abstract syntactic structure, before passing the result on to the cognitive-deductive system. In SHRDLU all three systems operate concurrently and are able to call upon each other at specific points.

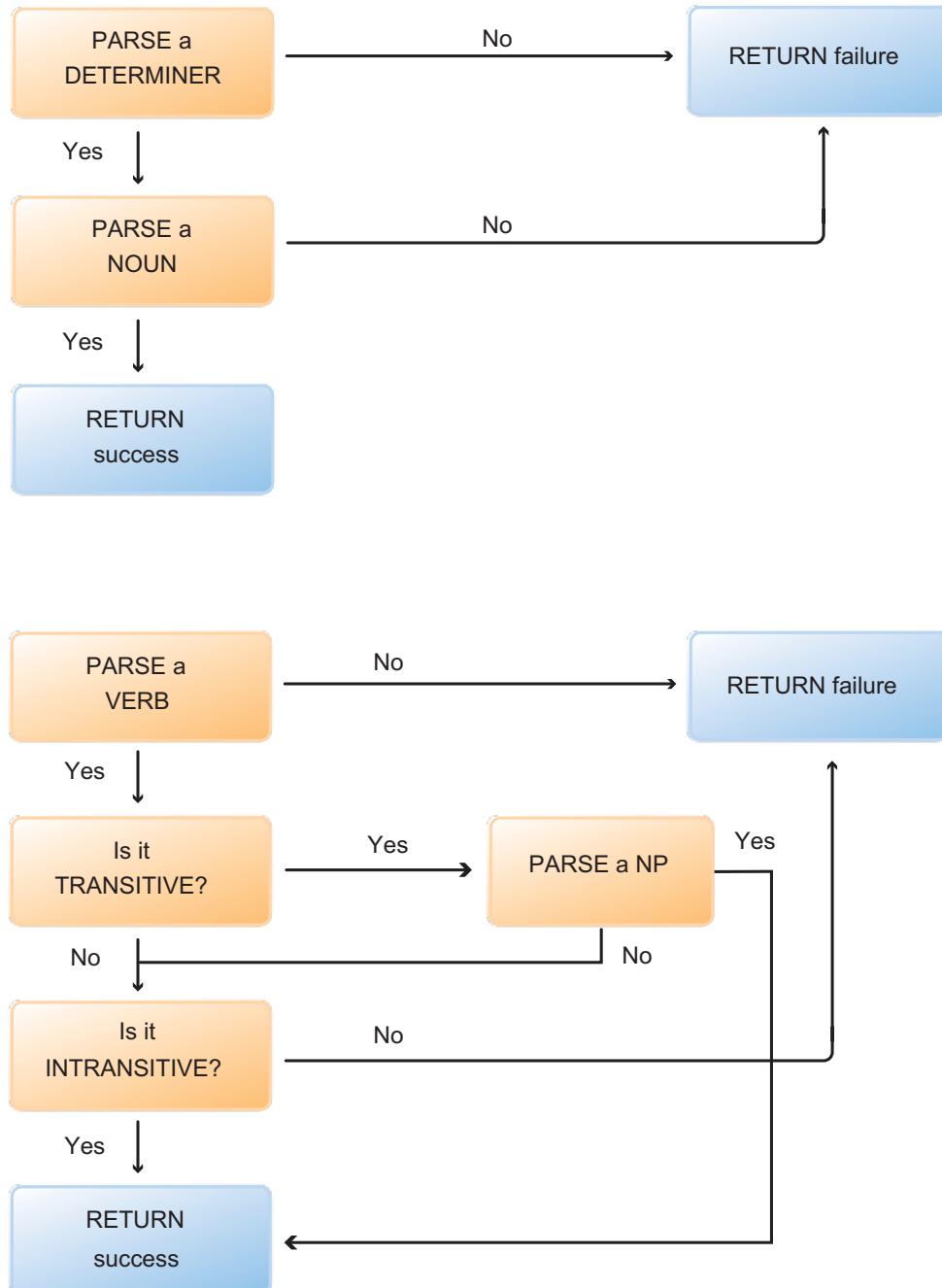
What makes this possible is that, although all three systems store and deploy different forms of knowledge, these different forms of knowledge are all represented in a similar way. They are all represented in terms of *procedures*.

The best way to understand what procedures are is to look at some examples. Let us start with the syntactic system, since this drives the whole process of language understanding. One very fundamental “decision” that the syntactic system has to make is whether its input is a sentence or not. Let us assume that we are dealing with a very simple language that only contains words in the following syntactic categories: Noun (e.g., “block” or “table”), Intransitive Verb (e.g., “\_\_\_\_ is standing up”), Transitive Verb (e.g., “\_\_\_\_ is supporting \_\_\_\_”), Determiner (e.g., “the” or “a”).

Figure 2.2 presents a simple procedure for answering this question. Basically, what the SENTENCE program does is exploit the fact that every grammatical sentence must contain



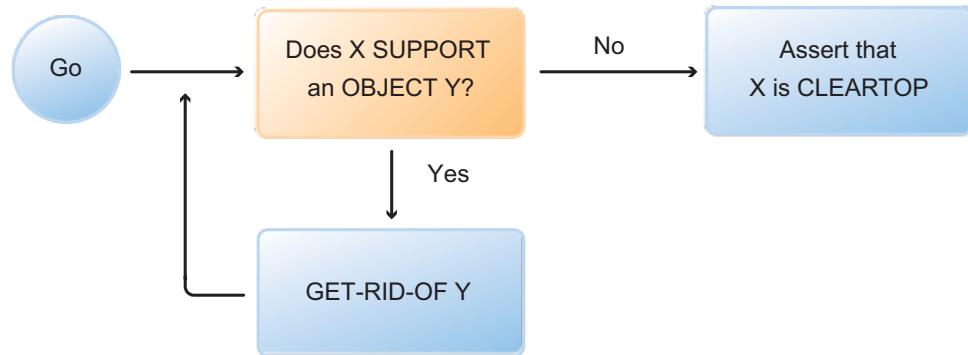
**Figure 2.2** An algorithm for determining whether a given input is a sentence or not. (Adapted from Winograd 1972)



**Figure 2.3** Algorithms for identifying noun phrases and verb phrases. (Adapted from Winograd 1973)

a noun phrase (NP) and a verb phrase (VP). It tests for the presence of a NP; tests for the presence of a VP; and then checks that there is no extra “junk” in the sentence.

Of course, in order to apply this procedure the syntactic system needs procedures for testing for the presence of noun phrases and verb phrases. This can be done in much the same way – by checking in an algorithmic manner whether the relevant syntactic units are present. Figure 2.3 gives two procedures that will work in our simple language.



**Figure 2.4** Procedure for applying the command CLEARTOP. (Adapted from Winograd 1972)

Moving to the job of semantic analysis, SHRDLU represents the meanings of words by means of comparable procedures. Instead of procedures for picking out syntactic categories, these procedures involve information about the micro-world and actions that the system can perform in the micro-world.

One of the words in SHRDLU's vocabulary is CLEARTOP. We can say that something (say, a block) is CLEARTOP when it does not have anything on it. CLEARTOP can also function as a command (as the command to remove anything resting on the block). CLEARTOP is represented in SHRDLU by the very simple procedure presented in Figure 2.4.

This definition of CLEARTOP exploits other “concepts,” such as SUPPORT and GET RID OF. Each of these other concepts has its own procedure, which may well call upon the CLEARTOP procedure.

To see how these procedures work to allow SHRDLU to follow instructions in the micro-world and answer questions about it we can look at the first few exchanges in a sample dialog described by Winograd in his 1973 paper. SHRDLU's contributions are in capital letters, while the sentences in italics were typed by a person. The commentary in normal type is by Winograd.

1. *Pick up a big, red block.*

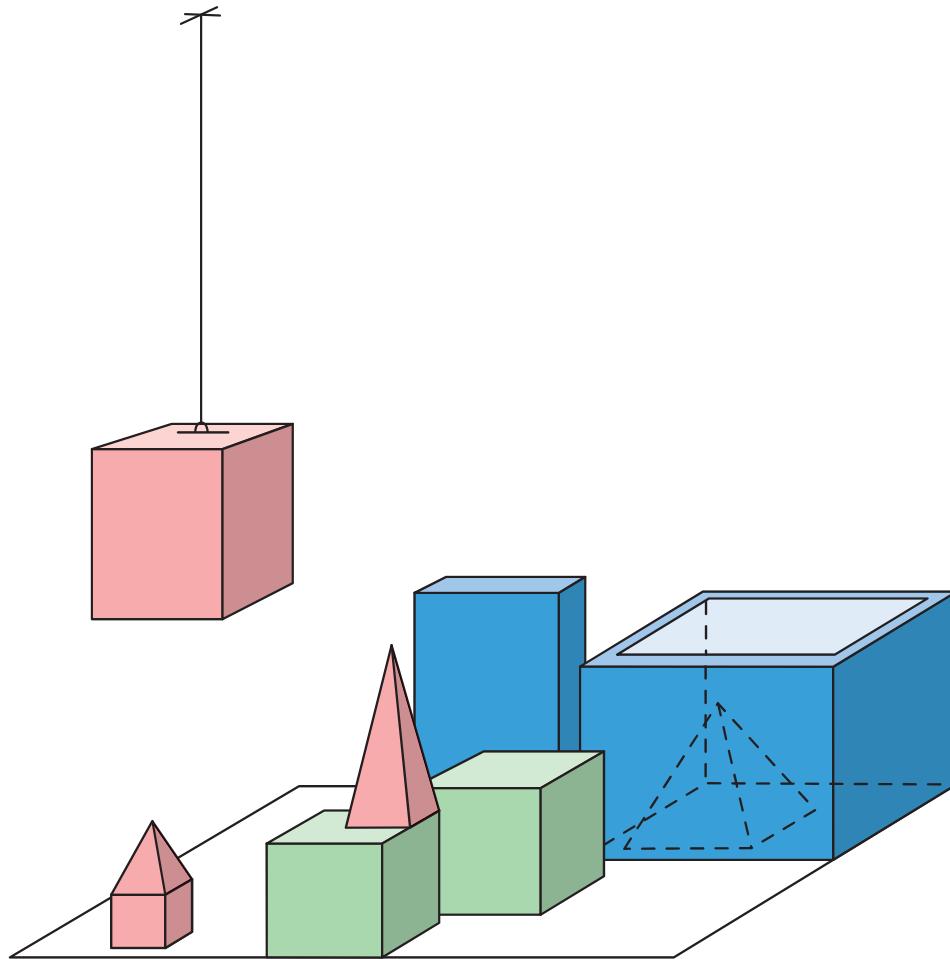
OK

The system answers “OK” when it carries out a command. In order to pick up the red block, it had to clear it off by finding a space for the green one and moving it away. The result is shown in Figure 2.5.

2. *Grasp the pyramid.*

I DON'T UNDERSTAND WHICH PYRAMID YOU MEAN.

There are three pyramids in the scene and the program knows that phrases beginning with “the” are intended to refer to a specific object the speaker has in mind. So SHRDLU does nothing.



**Figure 2.5** SHRDLU acting on the initial command to pick up a big red block. See the dialog in the text for what led up to this. (Adapted from Winograd 1972: 8)

3. *Find a block which is taller than the one you are holding and put it in the box.*

The situation is depicted in Figure 2.5. The system has to figure out what is meant by "one" and "it," as well as use deduction to decide which block it is holding.

BY "IT" I ASSUME YOU MEAN THE BLOCK WHICH IS TALLER THAN THE ONE I AM HOLDING

There were two possible referents for "it": the one selected, and "*the one you are holding*,"

OK

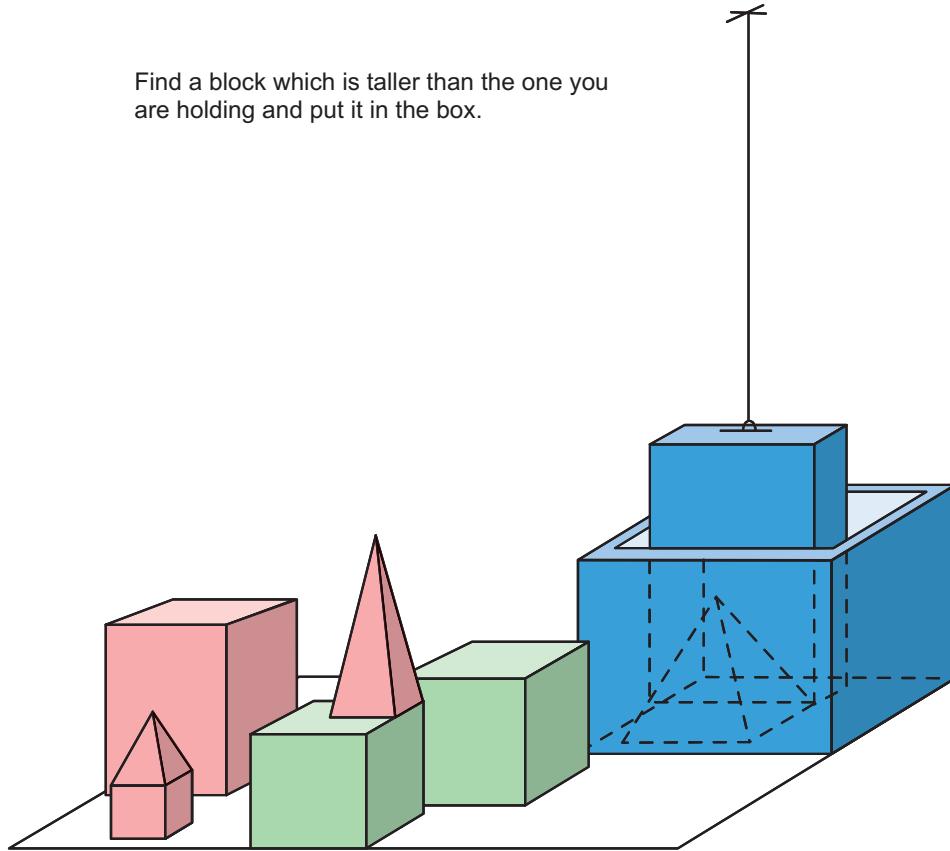
The hand has to set down the red block to pick up the taller blue one. We are now at Figure 2.6.

4. *What does the box contain?*

THE BLUE PYRAMID AND THE BLUE BLOCK

5. *What is the pyramid supported by?*

THE BOX



**Figure 2.6** SHRDLU completing instruction 3 in the dialog: “Find a block which is taller than the one you are holding and put it in the box.” (Adapted from Winograd 1972: Figure 3)

Earlier the system responded to “*the pyramid*” by saying “I DON’T UNDERSTAND WHICH PYRAMID YOU MEAN.” Now SHRDLU has identified the pyramid that was meant – the pyramid just mentioned in the response to the previous question.

Certainly, SHRDLU has its limitations. No one (certainly not Winograd) has ever claimed that it provides an accurate model of human language use – or even of a small segment of human language use. As Winograd himself was quick to acknowledge, SHRDLU does not really do justice to how conversations actually work in real life. Conversations are social interactions, not simply sequences of unconnected questions and answers. They involve agents trying to make sense of the world and of each other simultaneously.

Every level of linguistic understanding involves assumptions and guesses about what the other partner in the conversation is trying to communicate. It also involves assumptions and guesses about what they are trying to achieve. These are not always the same. In making and assessing those assumptions and guesses we use all sorts of heuristics and principles. We tend to assume, for example, that people generally tell the truth; that they don’t say things that are pointless and uninformative; and that what they say reflects what they are doing more generally. This is all part of what linguists call the *pragmatics* of



conversation. But there is nothing in SHRDLU's programming that even attempts to do justice to pragmatics.



### Exercise 2.2 Can you identify other ways in which SHRDLU falls short of modeling human understanding of natural language?

But to criticize SHRDLU for neglecting pragmatics, or for steering clear of complex linguistic constructions such as counterfactuals (statements about what *would* have happened, had things been different) is to miss what is genuinely pathbreaking about it.

SHRDLU illustrates a view of linguistic understanding as resulting from the interaction of many, independently specifiable cognitive processes. Each cognitive process does a particular job – the job of identifying noun phrases, for example. We make sense of the complex process of understanding a sentence by seeing how it is performed by the interaction of many simpler processes (or procedures). These cognitive processes are themselves understood algorithmically (although this is not something that Winograd himself stresses). They involve processing inputs according to rules. Winograd's procedures are sets of instructions that can be followed mechanically, just as in the classical model of computation (see Section 1.2).



## 2.2 How Do Mental Images Represent?

One way to understand a complex cognitive ability is to try to build a machine that has that ability (or at least some primitive form of it). The program that the machine runs is a model of the ability. Often the ability being modeled is a very primitive and simplified form of the ability that we are trying to understand. This is the case with SHRDLU, which was intended to model only a very basic form of linguistic understanding. But even in cases like that, we can still learn much about the basic principles of cognitive information processing by looking to see how well the model works. This is why the history of cognitive science has been closely bound up with the history of artificial intelligence.

We can think of artificial intelligence, or at least some parts of it, as a form of experimentation. Particular ideas about how the mind works are written into programs and then we "test" those ideas by seeing how well the programs work. But artificial intelligence is not the only way of developing and testing hypotheses open to cognitive scientists. Cognitive scientists have also learned much from the much more direct forms of experiment carried out by cognitive psychologists.

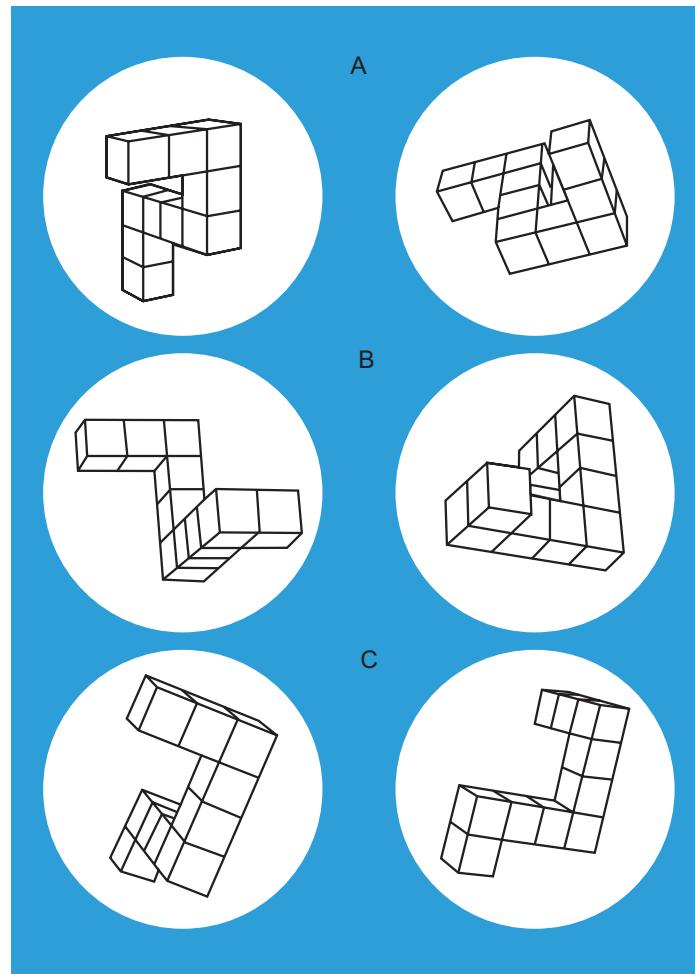
As we saw in the previous chapter, the emergence of cognitive psychology as a serious alternative to behaviorism in psychology was one of the key elements in the emergence of cognitive science. A good example of how cognitive psychology can serve both as an inspiration and as a tool for cognitive science came with what has come to be known as the imagery debate.

The imagery debate began in the early 1970s, inspired by a thought-provoking set of experiments on mental rotation carried out by the psychologist Roger Shepard in collaboration with Jacqueline Metzler, Lynn Cooper, and other scientists. The initial experiments (and many of the follow-up experiments) are rightly recognized as classics of cognitive

psychology. From the perspective of cognitive science, however, what is most interesting about them is the theorizing to which they gave rise about the format in which information is stored and the way in which it is processed. This was one of the first occasions when cognitive scientists got seriously to grips with the nature and format of mental representation – a theme that has dominated cognitive science ever since.

## Mental Rotation: Shepard and Metzler, "Mental Rotation of Three-Dimensional Objects" (1971)

The original mental rotation experiments are easy to describe. Subjects were presented with drawings of pairs of three-dimensional figures. Figure 2.7 contains examples of these pairs.



**Figure 2.7** Examples of the three-dimensional figures used in Shepard and Metzler's 1971 studies of mental rotation. Subjects were asked to identify which pairs depicted the same figure at different degrees of rotation. (Adapted from Shepard and Metzler 1971)



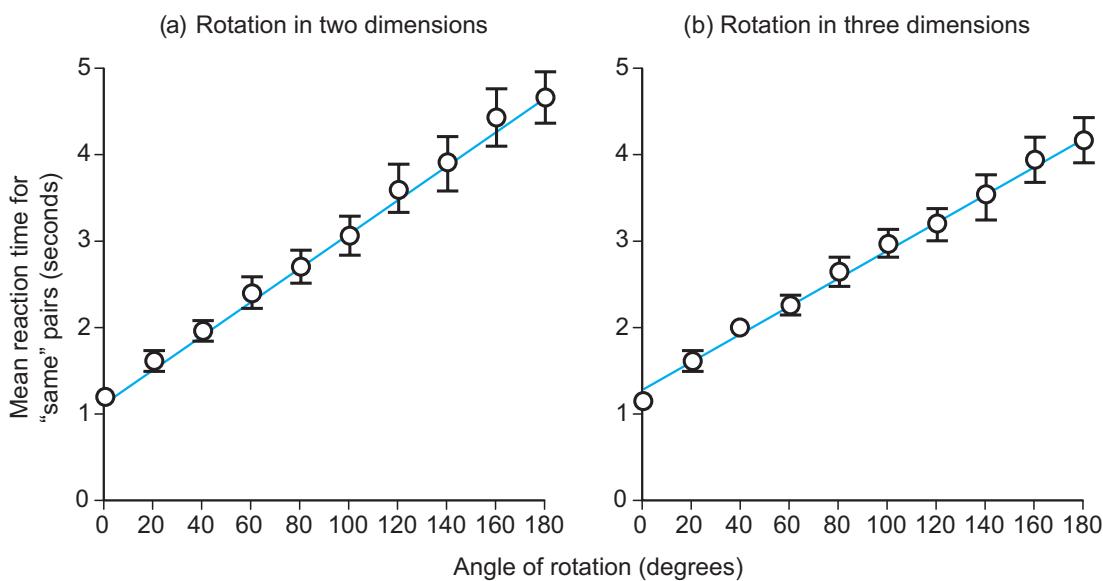
Each figure is asymmetric and resembles its partner. In two cases the figures resemble each other because they are in fact the same figure at different degrees of rotation. In a third case the figures are different. The subjects were asked to identify as quickly as possible pairs of drawings where the second figure is the same as the first, but rotated to a different angle. (You can do this experiment for yourself. Several versions of the Shepard–Metzler paradigm can be carried out online. See the Further Reading for an example. Putting “mental rotation” into a search engine will find others.)



### **Exercise 2.3 Which pair is the odd one out? In the pair with two distinct figures, how are those figures related to each other?**

Shepard and Metzler found that there is a direct, linear relationship between the length of time that subjects took to solve the problem and the degree of rotation between the two figures (see Figure 2.8). The larger the angle of rotation (i.e., the further the figures were from each other in rotational terms), the longer subjects took correctly to work out that the two drawings depicted the same figure. And the length of time increased in direct proportion to the degree of rotation. These findings have proved very robust. Comparable effects have been found in many follow-up experiments. Much more controversial is how to interpret what is going on.

The subjects in the original experiment were not asked to solve the problem in any particular way. They were simply asked to pull one lever if the two pictures represented the same figure, and another lever if the pictures represented different figures. The explanation that comes quickest to mind, though, is that the subjects solved the problem by mentally rotating one figure to see whether or not it could be mapped onto the other. This would



**Figure 2.8** Results of Shepard and Metzler’s 1971 studies of mental rotation. (a) Mean reaction time for shape rotation in two dimensions. (b) Mean reaction time for shape rotation in three dimensions.

certainly provide a neat explanation of the findings. And this is indeed how Shepard, Metzler, and many others did interpret them (not least because that is what many of the subjects described themselves as doing). This interpretation of the experiments raises some fundamental questions about the format in which information is encoded and manipulated in tasks of this type.



### **Exercise 2.4** Present in your own words Shepard and Metzler's conclusion. Explain their reasoning. What sort of assumptions does it rest on?

Suppose that we take the subject's report of what they are doing in the experiments at face value. Suppose, that is, that we think of the subjects as rotating mental images in their "mind's eye." It seems on the face of it that this is really just an application of a skill that we use all the time – the skill of transforming mental images in order to calculate, for example, whether one's car will fit into a tight parking space, or where a tennis ball will land. The question is not really whether we have such skills and abilities, but rather what makes them possible. And this is really a question about how the brain processes information.

The rotation in my "mind's eye" does not explain how I solve the problem. It is itself something that needs to be explained. What is the cognitive machinery that makes it possible for me to do what I might describe to someone else as rotating the mental image of a shape? Most cognitive scientists think that our conscious experience of rotating a mental image is the result of unconscious information processing. Information about the shape is derived from perception and then transformed in various ways that enable the subject to determine whether the two drawings are indeed drawings of the same shape. But the question is: How is that information represented and how is it transformed?



## Information Processing in Mental Imagery

The standard way of thinking about the mind as an information processor takes the digital computer as a model. (This was almost unchallenged in the early 1970s, and remains a popular view now, although we now have a much clearer sense of some alternative ways of thinking about information processing.) Digital computers store and manipulate information in a fixed format. Essentially, all forms of information in a digital computer are represented using the binary numerals 0 and 1.

Each binary digit carries a single unit of information (a *bit*). Within the computer these units of information are grouped into words – a *byte*, for example, is an 8-bit word that can carry 256 units of information. This way of carrying information in discrete quantities is often called digital information storage.

One feature of digitally encoded information is that the length of time it takes to process a piece of information is typically a function only of the quantity of information (the number of bits that are required to encode it). The particular information that is encoded ought not to matter. But what the mental rotation experiments seem to show is that there are information-processing tasks that take varying amounts of time even though the quantity of information remains the same.



### Exercise 2.5 Why does a byte carry 256 units of information? (Hint: There are eight digits, each of which can be in one of two states – so how many possible states are there?)

In order to get an intuitive picture of what is going on here and why it might seem puzzling, look again at the experimental drawings in Figure 2.7 and think about how each of them might be digitally encoded. Suppose that we think of each drawing as divided into many small boxes (rather like pixels on a television screen or computer monitor). Since the drawings are in black and white we can convey a lot of information about the drawing by stating, for each pixel, whether it is black or white. But this will not give us a full characterization, since the figures are represented three-dimensionally. This means that our characterization of each pixel that represents part of a surface will have to include a value for the surface's degree of orientation, degree of brightness, and so on.

Now, suppose that this has been done and that we have a pixel-by-pixel description of each drawing. This will be a collection of pixel descriptions. Each pixel description is simply a set of numbers that specifies the values on the relevant dimensions at the particular pixel locations. The overall pixel-by-pixel description of each drawing puts all those individual descriptions into an ordering that will allow it to be mathematically manipulated. One way of doing this would be to assign a set of coordinates to each pixel. In any event, the point is that each drawing can be represented by a set of numbers.

If this is how information is encoded, then solving the problem is essentially a matter of comparing two numerical descriptions to see if one can be mapped onto the other. Solving this problem is a tricky piece of mathematics that we fortunately do not have to go into, but there is no obvious reason why it should take longer to solve the problem for pairs of figures that are at greater degrees of rotation from each other than for pairs that are at smaller degrees from each other – and certainly no reason why there should be a linear relationship between reaction time and degree of rotation.

For reasons such as these, then, many cognitive scientists have suggested that mental rotation tasks tap into ways of encoding information very different from how information is encoded in a digital computer.

One distinctive feature of how information is represented in digital computers (what is often called digital representation) is that the connection between what we might think of as the unit of representation and what that unit represents is completely arbitrary. There is no reason, for example, why we should use the symbol “0” to represent a black pixel and the symbol “1” to represent a white pixel, rather than the other way around. The symbol “0” represents a black pixel because that is how the computer has been set up.

Contrast this with how, for example, a map represents a geographical region. Here there is a large-scale resemblance between the principal geographical features of the region and the discernible features of the map – if there is no such resemblance then the map will not be much use. The weaving and winding of a river is matched by the weaving and winding of the line on the map that represents the river. The outlines of a region of forestry are matched by the edges of the green patch on the map. Undulations in the terrain can be mapped onto the contour lines. And so on. A map is an excellent example of an imagistic

representation. The basic characteristic of an imagistic representation is that representation is secured through resemblance.

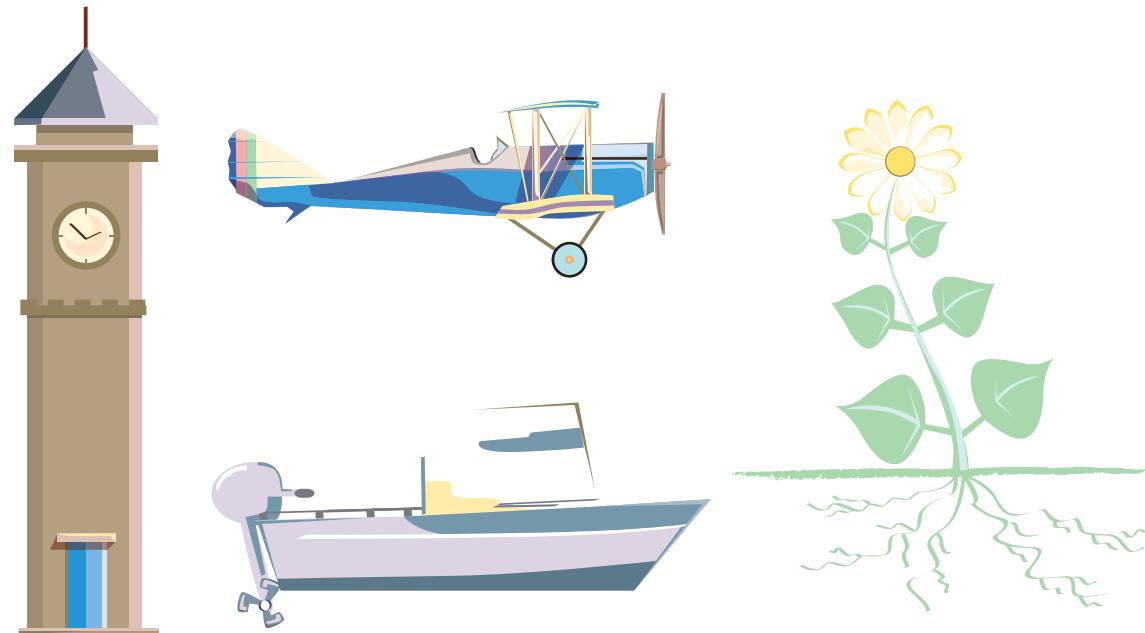


### **Exercise 2.6** Can you think of other differences between digital representation and imagistic representation?

One popular interpretation of the mental rotation experiments is as showing that at least some types of information are represented imagistically at the level of subconscious information processing. It is not just that we have the experience of consciously rotating figures in our mind's eye. The shapes are also represented imagistically in the subconscious information processing that makes possible these types of conscious experience.

The point of this interpretation is that certain operations can be carried out on imagistically represented information that cannot be carried out on digitally represented information. So, for example, it is relatively straightforward to think of rotating an imagistic representation, but as we saw earlier, difficult to think of rotating a digital representation. This gives us one way of explaining what is going on in the mental rotation experiments.

The idea that the information processing in mental imagery involves operations on imagistic representations also makes sense of many of the other effects identified in the experimental literature on imagery. So, for example, in a famous experiment carried out by Stephen Kosslyn in 1973 subjects were asked to memorize a set of drawings like those illustrated in Figure 2.9.



**Figure 2.9** Examples of vertically and horizontally oriented objects that subjects were asked to visualize in Kosslyn's 1973 scanning study. (Adapted from Kosslyn, Thompson, and Ganis 2006)



Kosslyn then gave them the name of one of the objects (e.g., “airplane”) and asked them to focus on one end of the memorized drawing. The experiment consisted of giving the subjects the names of possible parts of the object (e.g., “propeller”) and asking them to examine their images to see whether the object drawn did indeed have the relevant part (which it did on 50 percent of the trials). The subjects pushed a button only if they did indeed see the named part in their image of the drawn object.

Kosslyn found an effect rather similar to that in the mental rotation studies, namely, that the length of time it took the subjects to answer varied according to the distance of the parts from the point of focus. If the subjects were asked to focus on the tail of the plane, it would take longer for them to confirm that the plane had a propeller than that there was not a pilot in the cockpit.

Kosslyn’s interpretation of his own experiment was that the type of information processing involved in answering the test questions involves scanning imagistic representations. Instead of searching for the answer within a digitally encoded database of information about the figures, the subjects scan an imaginistically encoded mental image of the airplane.



### **Exercise 2.7** Can you think of a way of explaining the results of Kosslyn’s experiments without the hypothesis of imaginistically encoded information?

This takes us to the heart of a fundamental issue in cognitive science. Almost all cognitive scientists agree that cognition is information processing. But the imagery debate shows that there are competing models of how information is stored and how it is processed. We will return to these issues in later chapters.



## **2.3 An Interdisciplinary Model of Vision**

The mind can be studied at many different levels. We can study the mind from the bottom up, beginning with individual neurons and populations of neurons, or perhaps even lower down, with molecular pathways whose activities generate action potentials in individual neurons, and then trying to build up from that to higher cognitive functions. Or we can begin from the top down, starting out with general theories about the nature of thought and the nature of cognition and working downward to investigate how corresponding mechanisms might be instantiated in the brain. On either approach one will proceed via distinct levels of explanation that often have separate disciplines corresponding to them. A fundamental problem for cognitive science is working out how to combine and integrate different levels of explanation.

### **Levels of Explanation: Marr’s Vision (1982)**

The earliest systematic approach to tackling this problem is David Marr’s model of the human visual system in his 1982 book *Vision: A Computational Investigation into the Human*

*Representation and Processing of Visual Information.* Marr's conception of how different levels of explanation connect up with each other has been deeply influential, as a blueprint for practicing scientists and as a model for understanding the nature of explanation in cognitive science.

Marr distinguishes three different levels for analyzing cognitive systems. At the top is the *computational level*. Here cognitive scientists analyze in very general terms the particular type of task that the system performs. The tasks of an analysis at the computational level are:

- 1 to translate a general description of the cognitive system into a specific account of the particular information-processing problem that the system is configured to solve, and
- 2 to identify the constraints that hold upon any solution to that information-processing task.

The guiding assumption here is that cognition is ultimately a matter of information processing. A computational analysis identifies the information with which the cognitive system has to begin (the *input* to that system) and the information with which it needs to end up (the *output* from that system).



### **Exercise 2.8** Think of a specific cognitive system and explain what it does in information-processing terms.

Marr calls the next level down the *algorithmic level*. The algorithmic level tells us how the cognitive system actually solves the specific information-processing task identified at the computational level. It tells us how the input information is transformed into the output information. It does this by giving algorithms that effect that transformation.

So, an algorithmic-level explanation takes the form of specifying detailed sets of information-processing instructions that will explain how, for example, information from the sensory systems about the distribution of light in the visual field is transformed into a representation of the three-dimensional environment around the perceiver.

In contrast, the principal task at the *implementational level* is to find a physical realization for the algorithm – that is to say, to identify physical structures that will realize the representational states over which the algorithm is defined and to find mechanisms at the neural level that can properly be described as computing the algorithm in question.



### **Exercise 2.9** Explain in your own words the difference between algorithmic and implementational explanations.

Table 2.1 is a table from Marr's book illustrating how the different levels of explanation fit together. Marr's approach is a classic example of *top-down* analysis. He starts with high-level analysis of the specific information-processing problems that the visual system confronts, as well as the constraints under which the visual system operates. At each stage of the analysis these problems become more circumscribed and more determinate. The suggestions offered at the algorithmic and implementational levels are motivated by discussions of constraint and function at the computational level – that is, by considering



**TABLE 2.1** A table illustrating the three different levels that Marr identified for explaining information-processing systems

COMPUTATIONAL THEORY	REPRESENTATION AND ALGORITHM	HARDWARE IMPLEMENTATION
<i>What is the goal of the computation, why is it appropriate, and what is the logic of the strategy by which it can be carried out?</i>	<i>How can this computational theory be implemented? In particular, what is the representation for the input and output, and what is the algorithm for the transformation?</i>	<i>How can the representation and algorithm be realized physically?</i>

*Note.* Each level has its own characteristic questions and problems. (From Marr 1982)

which features of the environment the organism needs to model and the resources it has available to it.

## Applying Top-Down Analysis to the Visual System

We can get a better sense of how this general model of top-down analysis works in practice by looking at how Marr applied it in thinking about human vision.

The first point to note is that Marr's model is very interdisciplinary. His thinking at the computational level about what the visual system does was strongly influenced by research into brain-damaged patients carried out by clinical neuropsychologists. In his book he explicitly refers to Elizabeth Warrington's work on patients with damage to the left and right parietal cortex – areas of the brain that when damaged tend to produce problems in perceptual recognition.

Warrington noticed that the perceptual deficits of the two classes of patient are fundamentally different. Patients with right parietal lesions are able to recognize and verbally identify familiar objects *provided that they can see them from familiar or "conventional" perspectives*. From unconventional perspectives, however, these patients would not only fail to identify familiar objects but would also vehemently deny that the shapes they perceived could possibly correspond to the objects that they in fact were. Figure 2.10 provides an example of conventional and unconventional perspectives.

Patients with left parietal lesions showed a diametrically opposed pattern. Although left parietal lesions are often accompanied by language problems, patients with such lesions tend to be capable of identifying the shape of objects. They are as successful as normal subjects on matching tasks, and are perfectly able to match conventional and unconventional representations of the same object.



**Figure 2.10** Two images of a bucket. A familiar/conventional view is on the left, and an unfamiliar/unconventional view is on the right. (From Warrington and Taylor 1973)

Marr drew two conclusions about how the visual system functions from Warrington's neuropsychological observations. First, information about the shape of an object must be processed separately from information about what the object is for and what it is called. Second, the visual system can deliver a specification of the shape of an object even when the perceiver is unable to recognize the object.

Here is Marr describing how he used these neuropsychological data to work out the basic functional task that the visual system performs.

Elizabeth Warrington had put her finger on what was somehow the quintessential fact about human vision – that it tells us about shape and space and spatial arrangement. Here lay a way to formulate its purpose – building a description of the shapes and positions of things from images. Of course, that is by no means all that vision can do; it also tells us about the illumination and about the reflectances of the surfaces that make the shapes – their brightnesses and colors and visual textures – and about their motion. But these things seemed secondary; they could be hung off a theory in which the main job of vision was to derive a representation of shape.

(Marr 1982: 7)

So, at the computational level, the visual system's basic task is to construct a representation of the three-dimensional shape and spatial arrangement of an object in a form that will allow that object to be recognized. Since ease of recognition is correlated with the ability to extrapolate from the particular vantage point from which an object is viewed, Marr concluded that this representation of object shape should be on an object-centered rather than an egocentric frame of reference (where an egocentric frame of reference is one



centered on the viewer). This, in essence, is the theory that emerges at the computational level.



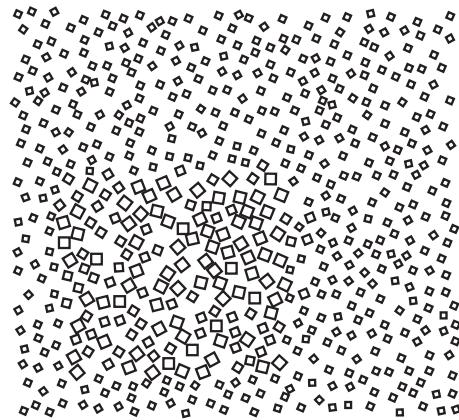
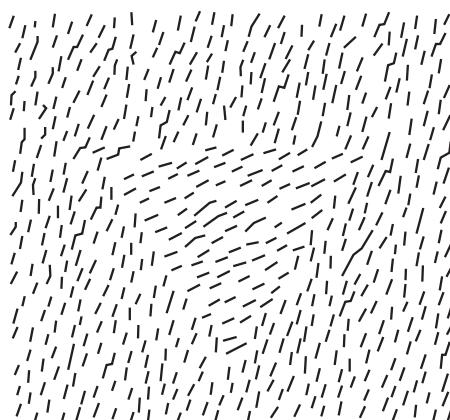
### Exercise 2.10 Explain in your own words why Marr drew the conclusions he did from Elizabeth Warrington's patients.

Analysis at the algorithmic level calls for a far more detailed account of how the general information-processing task identified at the computational level is carried out. What we are looking for now is an algorithm that can take the system from inputs of the appropriate type to outputs of the appropriate type. This raises a range of new questions. How exactly is the input and output information encoded? What are the system's *representational primitives* (the basic "units" over which computations are defined)? What sort of operations is the system performing on those representational primitives to carry out the information-processing task?

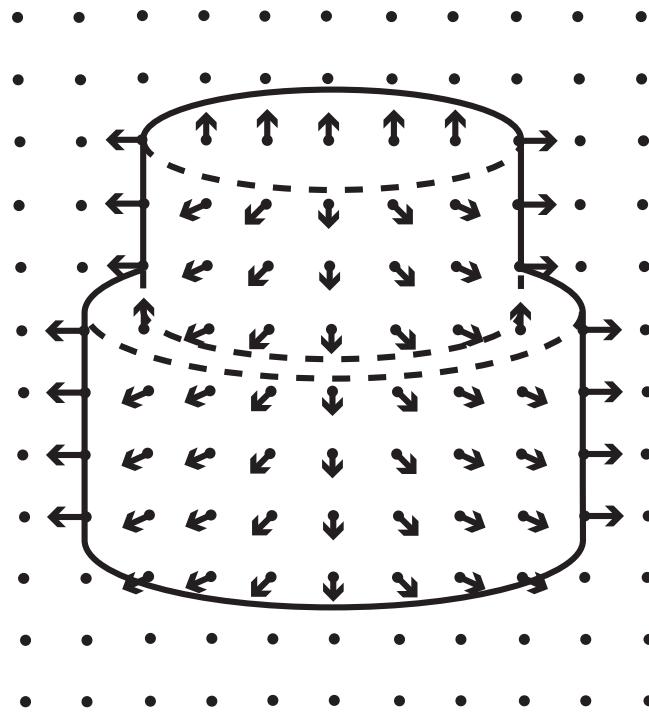
A crucial part of the function of vision is to recover information about surfaces in the field of view – in particular, information about their orientation; how far they are from the perceiver; and how they reflect light. In Marr's theory this information is derived from a series of increasingly complex and sophisticated representations, which he terms the *primal sketch*, the *2.5D sketch*, and the *3D sketch*.

The primal sketch makes explicit some basic types of information implicitly present in the retinal image. These include distributions of light intensity across the retinal image – areas of relative brightness or darkness, for example. The primal sketch also aims to represent the basic geometry of the field of view. Figure 2.11 gives two illustrations. Note how the primal sketch reveals basic geometrical structure – an embedded triangle in the left figure and an embedded square in the right.

The next information-processing task is to extract from the primal sketch information about the depth and orientation of visible surfaces from the viewer's perspective. The



**Figure 2.11** Two examples of Marr's primal sketch, the first computational stage in his analysis of the early visual system. The primal sketch contains basic elements of large-scale organization (the embedded triangle in the left-hand sketch, for example). (Adapted from Marr 1982)



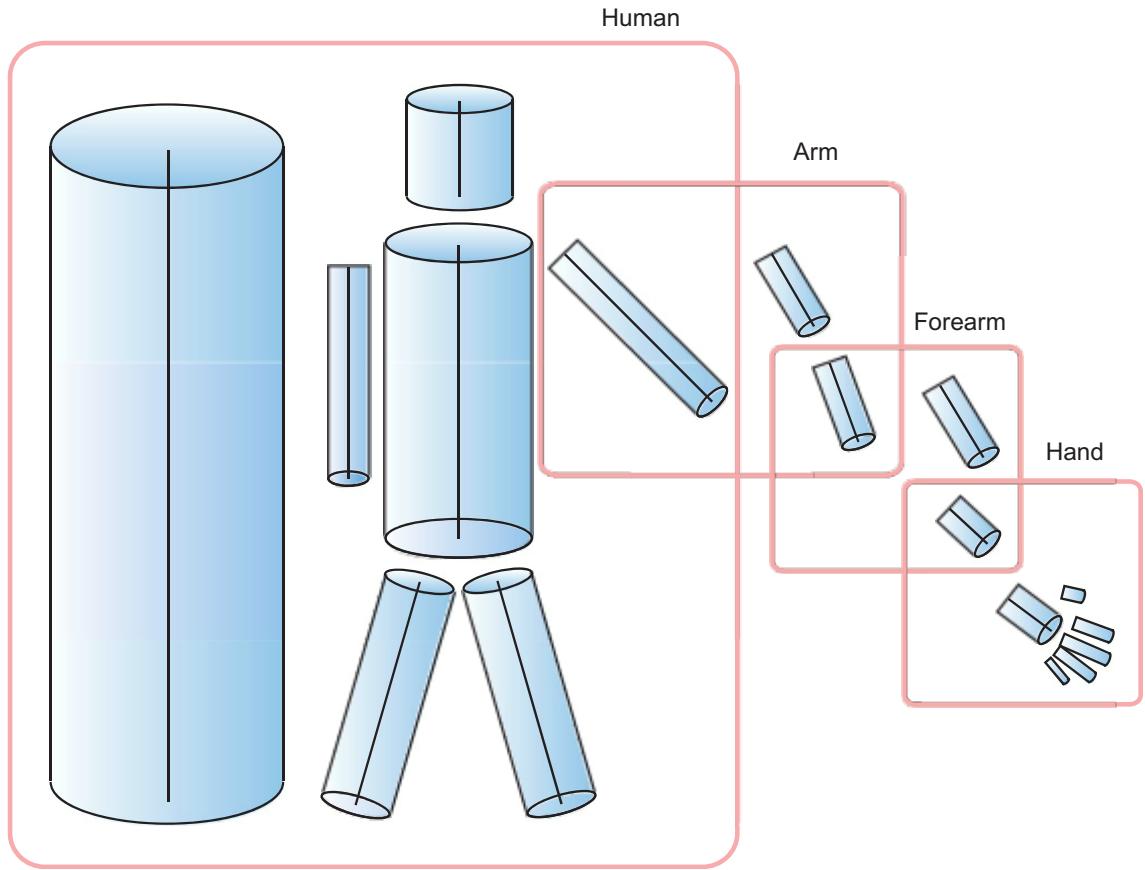
**Figure 2.12** An example of part of the 2.5D sketch. The figure shows orientation information but no depth information. (Adapted from Marr 1982)

result of this information processing is the 2.5D sketch. For every point in the field of view, the 2.5D sketch represents how far it is from the observer, as illustrated in Figure 2.12.

The 2.5D sketch is viewer-centered. It depends upon the viewer's particular vantage point. One of the crucial things that the visual system allows us to do, though, is to keep track of objects even though their visual appearance changes from the viewer's perspective (because either the object or the viewer is moving, for example). This requires a stable representation of object shape that is independent of the viewer's particular viewpoint. This viewer-independent representation is provided by the 3D sketch, as illustrated in Figure 2.13.

At the algorithmic level the job is to specify these different sketches and explain how the visual system gets from one to the next, starting with the basic information arriving at the retina. Since the retina is composed of cells that are sensitive to light, this basic information is information about the intensity of the light reaching each of those cells.

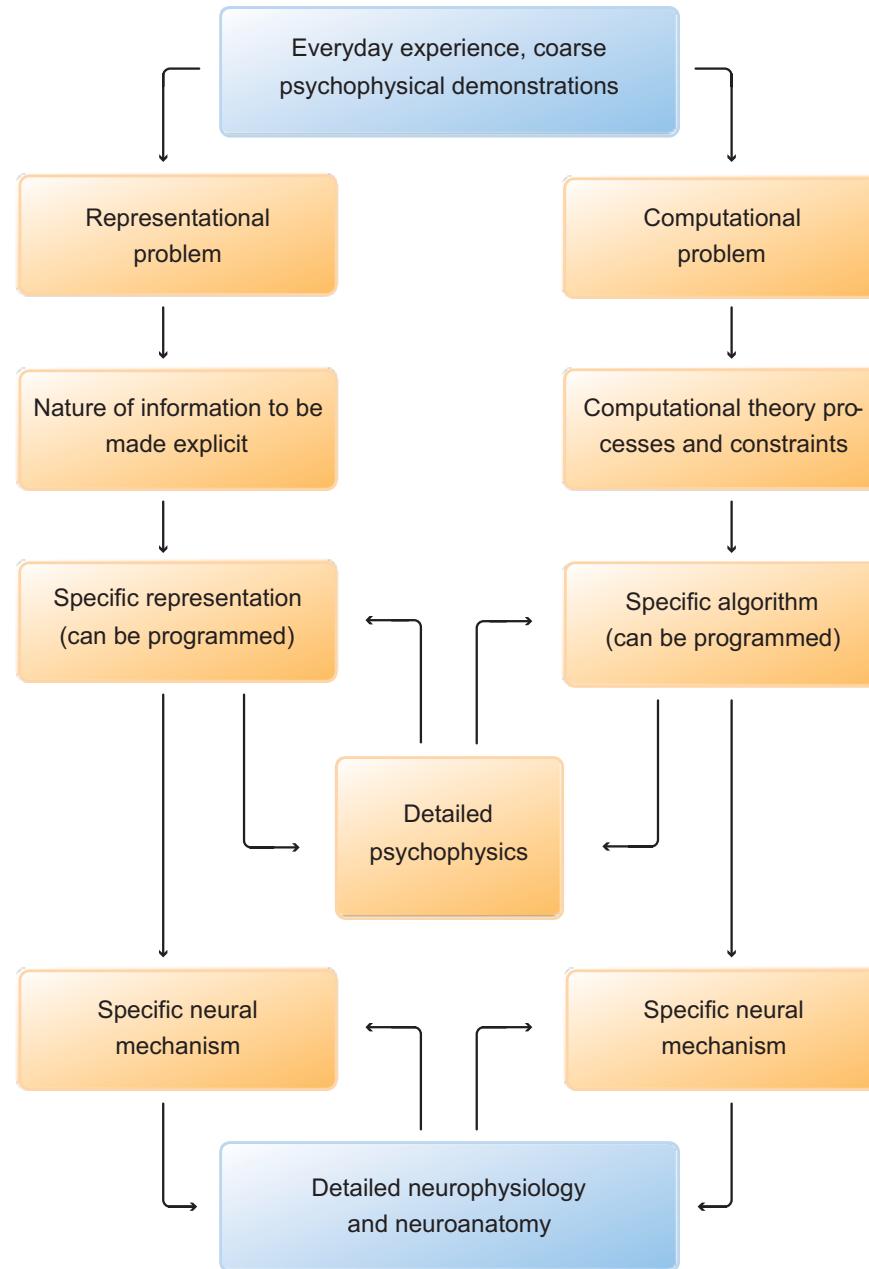
What are the starting-points for the information processing that will yield as its output an accurate representation of the layout of surfaces in the distal environment? Marr's answer is that the visual system needs to start with discontinuities in light intensity, because these are a good guide to boundaries between objects and other physically relevant properties. Accordingly the representational primitives that he identifies are all closely correlated with changes in light intensity.



**Figure 2.13** An illustration of Marr's 3D sketch, showing how the individual components are constructed. The 3D sketch gives an observer-independent representation of object shape and size. (Adapted from Marr 1982)

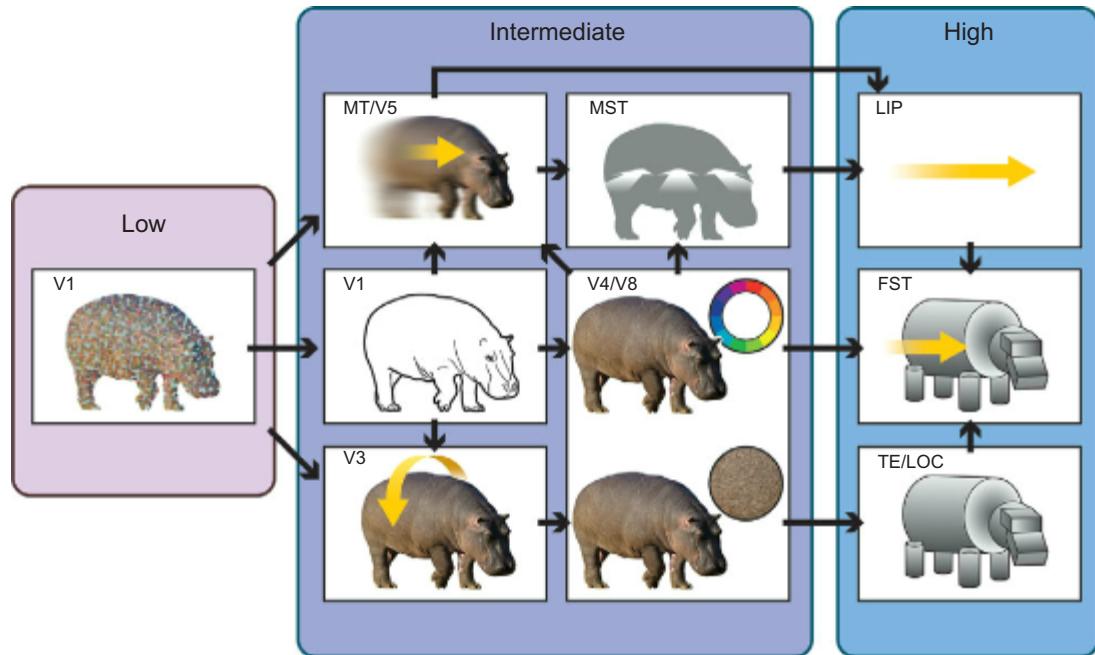
These representational primitives include *zero-crossings* (registers of sudden changes in light intensity), blobs, edges, segments, and boundaries. The algorithmic description of the visual system takes a representation formulated in terms of these representational primitives as the input, and spells out a series of computational steps that will transform this input into the desired output, which is a representation of the three-dimensional perceived environment.

Moving down to the implementational level, a further set of disciplines come into play. In thinking about the cognitive architecture within which the various algorithms computed by the visual system are embedded, we will obviously need to take into account the basic physiology of the visual system – and this in turn is something that we will need to think about at various different levels. Marr's own work on vision contains relatively little discussion of neural implementation. But the table from his book shown here as Figure 2.14 illustrates where the implementational level fits into the overall picture. Figure 2.15 is a more recent attempt at identifying the neural structures underlying the visual system.



**Figure 2.14** The place of the implementational level within Marr's overall theory. Note also the role he identifies for detailed experiments in psychophysics (the branch of psychology studying how perceptual systems react to different physical stimuli). (Adapted from Marr 1982)

Marr's analysis of the visual system clearly illustrates both how a single cognitive phenomenon can be studied at different levels of explanation, and how the different levels of explanation can come together to provide a unified analysis. It is not surprising that Marr's analysis of the visual system has been taken as a paradigm of how cognitive science ought to proceed.



#### Key:

- V1–V8: areas of the visual cortex in the occipital lobe (the back of the head). V1 produces the color and edges of the hippo but no depth. V2 produces the boundaries of the hippo. V3 produces depth. V4/V8 produces color and texture.
- MT: medial temporal area (often used interchangeably with V5). Responsible for representing motion.
- MST: medial superior temporal area. Responsible for representing size of the hippo as it gets nearer in space.
- LIP: lateral intraparietal area. Registers motion trajectories.
- FST: fundus of the superior temporal sulcus. Discerns shape from motion.
- TE: temporal area. Along with LOC, is responsible for shape recognition.
- LOC: lateral occipital complex

**Figure 2.15** An illustration of the hierarchical organization of the visual system, including which parts of the brain are likely responsible for processing different types of visual information. (From Prinz 2012)



## Summary

This chapter continued our historical overview of key steps in the emergence and evolution of cognitive science. We have reviewed three case studies: Terry Winograd's SHRDLU program for modeling natural language understanding; the explorations into the representational format of mental imagery inspired by the mental rotation experiments of Roger Shepard and others; and the multilevel analysis of the early visual system proposed by David Marr. Each of these represented a significant milestone in the emergence of cognitive science. In their very different ways they show

how researchers brought together some of the basic tools discussed in Chapter 1 and applied them to try to understand specific cognitive capacities.

## Checklist

### Winograd's SHRDLU

- (1) SHRDLU is more sophisticated than a conversation-simulating chatbot because it uses language to report on the environment and to plan action.
- (2) SHRDLU illustrated how abstract grammatical rules might be represented in a cognitive system and integrated with other types of information about the environment.
- (3) The design of SHRDLU illustrates a common strategy in cognitive science, namely, analyzing a complex system by breaking it down into distinct components, each performing a circumscribed information-processing task.
- (4) These information-processing tasks are implemented algorithmically (as illustrated by the flowcharts that Winograd used to explain SHRDLU's different procedures).

### The Imagery Debate

- (1) The experiments that gave rise to the imagery debate forced cognitive scientists to become much more reflective about how they understand information and information processing.
- (2) The imagery debate is not a debate about conscious experiences of mental imagery. It is about the information processing underlying those conscious experiences.
- (3) The mental rotation and scanning experiments were taken by many cognitive scientists to show that some information processing involves operations on geometrically encoded representations.
- (4) The debate is about whether the different effects revealed by experiments on mental imagery can or cannot be explained in terms of digital information-processing models.

### Marr's Theory of Vision

- (1) Marr identified three different levels for analyzing cognitive systems.
- (2) His analysis of vision is a classic example of the top-down analysis of a cognitive system. The analysis is driven by a general characterization at the computational level of the information-processing task that the system is carrying out.
- (3) This general analysis at the computational level is worked out in detail at the algorithmic level, where Marr explains how the information-processing task can be algorithmically carried out.
- (4) The bottom level of analysis explains how the algorithm is actually implemented. It is only at the implementational level that neurobiological considerations come directly into the picture.

## Further Reading

The general historical works mentioned at the end of the previous chapter also cover the material in this chapter and will provide further useful context-setting.

A web-based version of ELIZA can be found in the online resources. The principal resource for SHRDLU is Winograd's book *Understanding Natural Language* (1972). This is very detailed,



however, and a more accessible treatment can be found in his article "A procedural model of language understanding" (1973), which is reprinted in Cummins and Cummins 2000. One of the important descendants of the micro-world strategy exploited in SHRDLU was research into expert systems. A helpful introduction is the entry on expert systems in the *Macmillan Encyclopedia of Cognitive Science* (Medsker and Schulte 2003). The online *Encyclopedia of Cognitive Science* (Nadel 2005) also has an entry on SHRDLU.

Many of the most important original articles in the imagery debate are collected in Block 1981. The experiments described in the text were originally reported in Shepard and Metzler 1971, Kosslyn 1973, and Cooper and Shepard 1973. Demonstrations and further discussion of mental imagery can be found in the online resources. The imagery debate has received a good deal of attention from philosophers. Rollins 1989 and Tye 1991 are book-length studies. The *Stanford Encyclopedia of Philosophy* also has an entry on mental imagery at <http://plato.stanford.edu/entries/mental-imagery/mental-rotation.html>. Kosslyn, Thompson, and Ganis 2006 is a recent defense of geometric representation from one of the central figures in the debate. The best meta-analyses of mental imagery studies can be found in Voyer, Voyer, and Bryden 1995 and Zacks 2008.

Marr's book on vision (1982) has recently been reprinted (2010). Shimon Ullman's foreword in the new edition and Tomaso Poggio's afterword provide some background to Marr. Ullman discusses where the field has moved since Marr, while Poggio discusses Marr's contribution to computational neuroscience and how the field can benefit from looking back to Marr. The first chapter of Marr's book is reprinted in a number of places, including Bermúdez 2006 and Cummins and Cummins 2000. Marr's selected papers have also been published together (Vaina 1991). Dawson 1998 is a textbook on cognitive science that is structured entirely around Marr's tri-level hypothesis. Also see Tsotsos 2011. Chapter 2 of Prinz 2012 gives a general assessment of the accuracy of Marr's account, in light of current research on visual processing. Elizabeth Warrington's classic studies can be found in Warrington and Taylor 1973, 1978.





## CHAPTER THREE

# The Turn to the Brain

### OVERVIEW 65

- 3.1 Cognitive Systems as Functional Systems? 66**
- 3.2 The Anatomy of the Brain and the Primary Visual Pathway 68**  
The Two Visual Systems Hypothesis:  
Ungerleider and Mishkin, "Two Cortical Visual Systems" (1982) 70
- 3.3 Extending Computational Modeling to the Brain 76**  
A New Set of Algorithms: Rumelhart, McClelland, and the PDP Research Group, *Parallel Distributed Processing: Explorations in the Microstructure of Cognition* (1986) 77  
Pattern Recognition in Neural Networks:  
Gorman and Sejnowski, "Analysis of Hidden Units in a Layered Network Trained to Identify Sonar Targets" (1998) 78

### 3.4 Mapping the Stages of Lexical Processing 80

- Functional Neuroimaging with PET 81  
Petersen, Fox, Posner, and Mintun, "Positron Emission Tomographic Studies of the Cortical Anatomy of Single-Word Processing" (1988) 81

### 3.5 Studying Memory for Visual Events 84

- Functional Neuroimaging with fMRI 86  
Brewer, Zhao, Desmond, Glover, and Gabrieli, "Making Memories: Brain Activity That Predicts How Well Visual Experience Will Be Remembered" (1998) 87

### 3.6 The Neural Correlates of the BOLD Signal 90

- Logothetis, "The Underpinnings of the BOLD Functional Magnetic Resonance Imaging Signal" (2001) 91



## Overview

A striking feature of contemporary cognitive science, as compared with the 1970s for example, is the increasing centrality of the brain. This chapter reviews some landmarks in cognitive science's turn to the brain. There are several different strands here. One is the emergence of different techniques for studying the brain. These include brain studies and functional neuroimaging techniques. And then, distinct from these but no doubt related, is the development of neurally inspired computational models.

For both theoretical and practical reasons, neuroscience was fairly peripheral to cognitive sciences until the 1980s. We begin in Section 3.1 by looking at some of the theoretical reasons, particularly the influential idea that cognitive systems are functional systems, and so need to be studied in terms of their function – what they do and how they do it. Many cognitive scientists hold

that this type of functional analysis can be carried out without looking at details of neural implementation.

One early move away from this functional view came with the two visual systems hypothesis, originally proposed by the neuroscientists Leslie Ungerleider and Mortimer Mishkin. In Section 3.2 we look at how Ungerleider and Mishkin drew conclusions about the structure and organization of vision from data about the pathways in the brain that carry visual information. The direction of explanation is bottom-up, rather than top-down (as in Marr's framework, which we looked at in Section 2.3).

An important factor in the turn toward the brain was the development of ways of modeling cognitive abilities designed to reflect very general properties of brains. As described in Section 3.3, so-called connectionist networks, or artificial neural networks, involve large populations of neuron-like units. The individual units are not biologically plausible in any detailed sense. But the network as a whole behaves in ways that reflect certain high-level properties of brain functioning.

Until the 1980s, techniques for studying human brains while cognitive tasks were actually being carried out were relatively unsophisticated and not widely known among cognitive scientists. This changed with the emergence of functional neuroimaging in the 1980s, which provided a powerful tool for studying what goes on in the brain when subjects are actually performing different types of cognitive task.

In Section 3.4 we look at an early and very influential application of positron emission tomography (PET) scanning technology to the study of visual word processing. This study shows how functional neuroimaging can be used to generate information-processing models of how cognitive tasks are carried out – information-processing models that are derived, not from abstract task analysis, but rather from detailed study of neural activity.

Section 3.5 introduces functional magnetic resonance imaging (fMRI), which has superseded PET in many domains, and allows a different type of experimental design (known as event-related design). Event-related fMRI is introduced through studies of visual memory.

Finally, in Section 3.6, we turn to what fMRI actually measures, looking at groundbreaking experiments by Nikos Logothetis. These experiments use single electrode recordings of individual neurons to study the type of brain activity that is correlated with the BOLD (blood oxygen level dependent) signal that is directly measured by fMRI (functional magnetic resonance imaging).



## 3.1

## Cognitive Systems as Functional Systems?

First, though, some background to make clear the significance of the turn to the brain. Most of the models that we have looked at so far in our historical survey share certain very basic features. In particular, models here assume that information is transformed and transmitted in the brain in much the same way as information is transformed and transmitted in digital computers. We can study computer algorithms without thinking about the hardware and circuitry on which they run. And so, it is not surprising that these models typically abstract away from the details of neural machinery in thinking about the algorithms of cognition.



In fact, for many cognitive scientists it is not just that cognitive processes *can be studied* independently of the neural machinery on which they run. They *have to be studied* that way. This is because they think of cognitive systems as *functional* systems. The important point is, as the word suggests, that functional systems are to be understood primarily in terms of their function – what they do and how they do it. And, these cognitive scientists emphasize, this type of analysis can be given without going into details about the particular physical structure implementing that function.

An analogy will help. Consider a heart. What makes something a heart? The most important thing is what it does. Hearts are organs that pump blood around the body – in particular, they collect deoxygenated blood and pump it toward the lungs where it becomes reoxygenated. The actual physical structure of the heart is not particularly important. An artificial heart will do the job just as well (although not perhaps for as long) and so still counts as a heart. Crocodiles and humans have hearts with four chambers, while most reptiles have hearts with three chambers. What matters is the job the heart does, not how it does it. A gray whale's heart is no less a heart than a hummingbird's heart just because the first beats 9 times per minute while the second beats 1,200 times per minute. One way of putting this is to say that functional systems are *multiply realizable*. The heart function can be realized by multiple different physical structures.



### Exercise 3.1 Give another example of a multiply realizable system.

If cognitive systems are functional systems that are multiply realizable in the way that the heart is multiply realizable, then, the argument goes, it is a mistake to concentrate on the details of how the brain works. In fact, according to cognitive scientists opposed to looking at the brain, focusing on how the brain works is likely to lead to a misleading picture of how cognition works. It might lead us to take as essential to memory, say, things that are really just contingent properties of how our brains have evolved. We would be making the same mistake as if we were to conclude that hearts have to have four chambers because the human heart does, or if we decided that Microsoft Word has to run on a 2.33 GHz Intel Core 2 Duo processor just because that is the processor in my Apple Macintosh.

But other cognitive scientists think that abstracting away from neural machinery in studying the algorithms of cognition may not be a good idea. For one thing, cognitive activity needs to be coordinated with behavior and adjusted online in response to perceptual input. The control of action and responsiveness to the environment requires cognitive systems with a highly developed sense of timing. The right answer is no use if it comes at the wrong time. But how can we think about the speed and efficiency of the mind without taking into account the fact that it runs on the hardware of the brain?

Moreover, the mind is not a static phenomenon. Cognitive abilities and skills themselves evolve over time, developing out of more primitive abilities and giving rise to further cognitive abilities. Eventually they deteriorate and, for many of us, gradually fade out of existence. In some unfortunate cases they are drastically altered as a result of traumatic damage. This means that an account of the mind must be compatible with plausible accounts of how cognitive abilities emerge. It must be compatible with what we know

about how cognitive abilities deteriorate. It must be compatible with what we know about the relation between damage to the brain and cognitive impairment.

Cognitive abilities tend to *degrade gracefully*. As we get older reaction times increase, motor responses slow down, and recall starts to become more problematic. But these abilities do not (except as a result of trauma or disease) disappear suddenly. The deterioration is gradual, incremental, and usually imperceptible within small time frames. This type of graceful degradation is a function of how brains are wired, and of the biochemistry of individual neurons.

The same holds for how cognitive abilities emerge and develop. Brains learn the way they do because of how they are constructed – and in particular because of the patterns of connectivity existing at each level of neural organization (between neurons, populations of neurons, neural systems, neural columns, and so forth). It is plausible to expect our higher-level theories of cognitive abilities to be constrained by our understanding of the neural mechanisms of learning.



**Exercise 3.2** Can you think of other reasons for thinking that we should not theorize about cognition without theorizing about the brain?



## 3.2

## The Anatomy of the Brain and the Primary Visual Pathway

We turn now to the two visual systems hypothesis, as our first illustration of the turn to the brain. First, though, we need a little information about the large-scale anatomy of the brain.

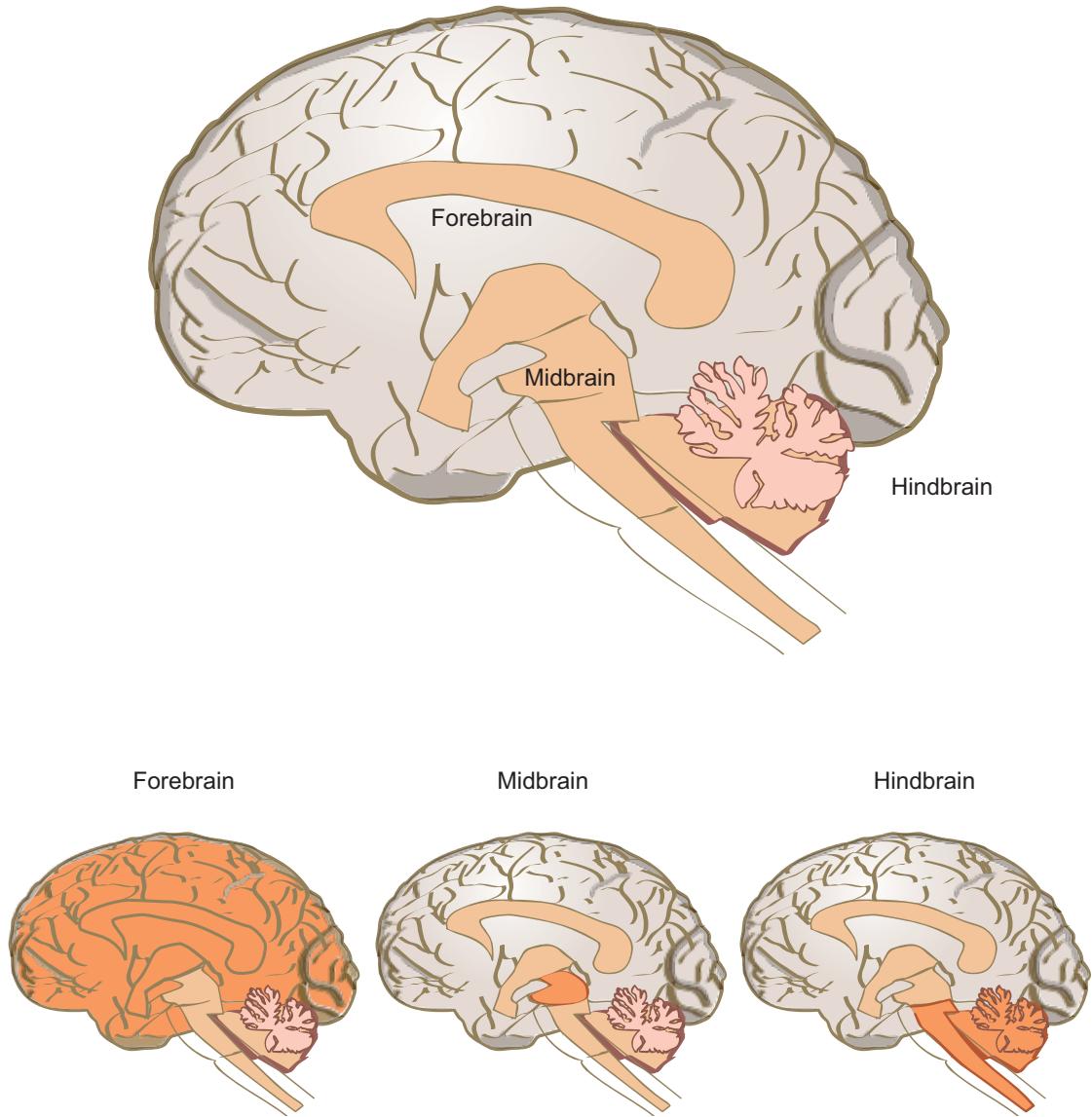
Anatomists distinguish three different parts of the mammalian brain – the *forebrain*, the *midbrain*, and the *hindbrain*. This structure is illustrated for the human brain in Figure 3.1.

As the figure shows, the forebrain is the largest of the three regions. Most of the forebrain is taken up by the *cerebrum* (see Figure 3.2), which is the main portion of the brain and the most important for cognitive and motor processing. The cerebrum is divided into two hemispheres – left and right. The hemispheres are separated by a deep groove, known as the *longitudinal fissure* or the *interhemispheric fissure*.

The two hemispheres have similar organizations. Each has an outer layer, which comprises what is known as the cerebral cortex. Only mammals have a cerebral cortex, and in the human brain it is about 2–4 mm thick. Moving inward from the outer, cortical layer we find the remaining major structures of the forebrain. These are the thalamus, the hypothalamus, and the limbic system, collectively known as subcortical areas (because they lie below the cortex).

The tissue of these inner, subcortical parts of the forebrain is known as white matter, because it is made up of neurons whose axons are encased in a white sheath of myelin (which speeds up the transmission of nerve signals). Neurons in the cerebral cortex are typically unmyelinated and look gray – hence the term “gray matter.”

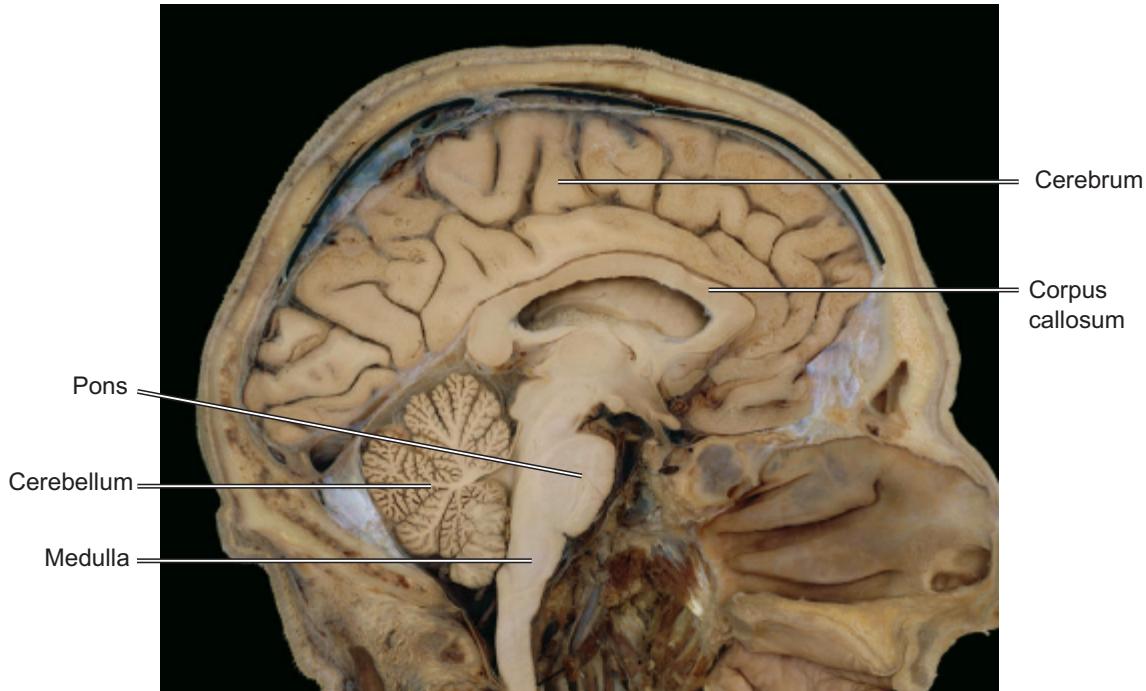
Within each hemisphere, the cerebral cortex is divided into four main regions, called *lobes*. Each lobe is believed to be responsible for carrying out different cognitive tasks. Figure 3.3 illustrates the organization of the left hemisphere into four lobes, while Box 3.1 summarizes what each lobe is believed to be specialized for.



**Figure 3.1** The large-scale anatomy of the brain, showing the forebrain, the midbrain, and the hindbrain.

### BOX 3.1 What Does Each Lobe Do?

- Frontal lobe – reasoning, planning, parts of speech, movement, emotions, and problem solving
- Parietal lobe – movement, orientation, recognition, perception of stimuli
- Occipital lobe – associated with visual processing
- Temporal lobe – associated with perception and recognition of auditory stimuli, memory, and speech



**Figure 3.2** A vertical slice of the human brain, showing the cerebrum. © TISSUEPIX/SCIENCE PHOTO LIBRARY

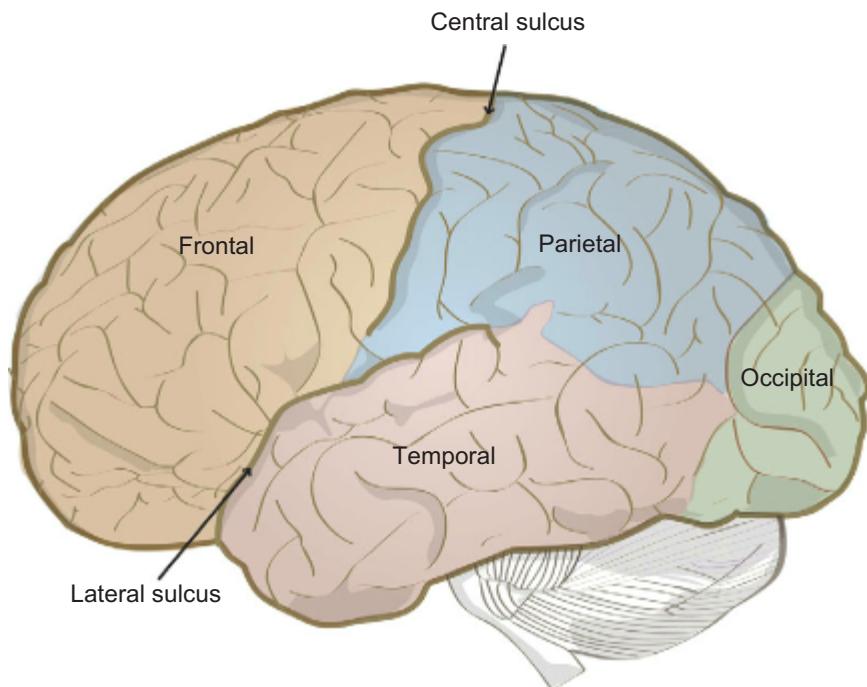
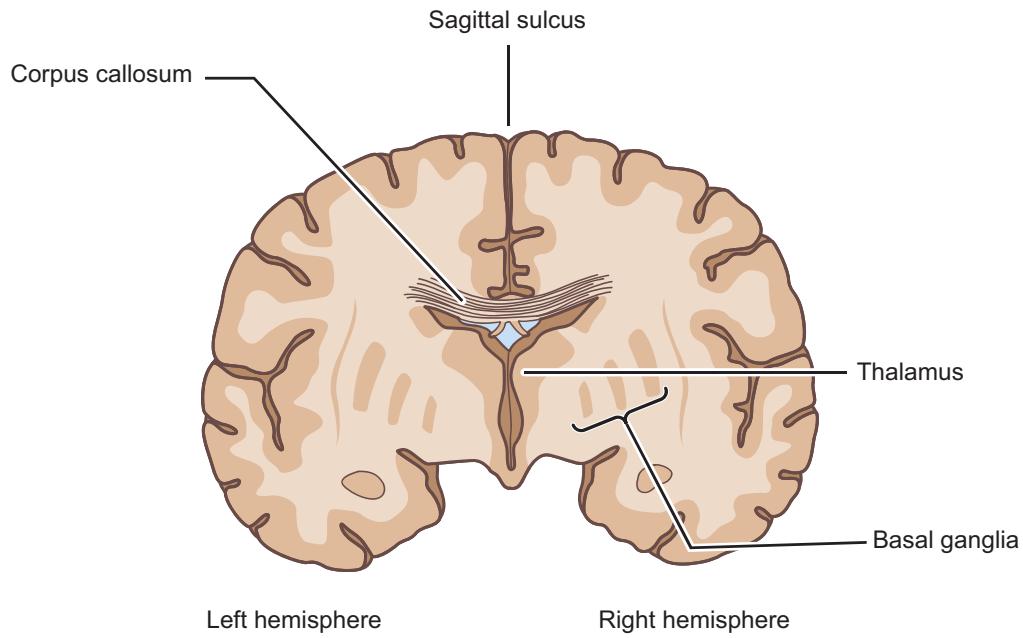
There is further organization within each lobe. In 1909 the German neurologist Korbinian Brodmann proposed a mapping of the cerebral cortex into fifty-two areas. These Brodmann areas are still in use today. An example particularly relevant to us now is Brodmann area 17, which is also known as the *primary visual cortex*, the *striate cortex*, or area V1. Brodmann area 17 is located in the *occipital lobe* and (as the name “primary visual cortex” suggests) it is the point of arrival in the cortex for information from the retina.

The information pathway leading from the retina to the primary visual cortex is relatively well understood. It is illustrated in Figure 3.4, which shows how visual information from each eye is transmitted by the optic nerve to the lateral geniculate nucleus (a subcortical area of the forebrain) and thence to the primary visual cortex. The diagram clearly shows the *contralateral* organization of the brain. Each hemisphere processes information deriving from the opposite side of space. So, visual information from the right half of the visual field is processed by the left hemisphere (irrespective of which eye it comes from).

Much more complicated than the question of how information from the retina gets to the primary visual cortex is the question of what happens to that information when it leaves the primary visual cortex. This is where we come to the two visual systems hypothesis and to the work of Ungerleider and Mishkin.

## The Two Visual Systems Hypothesis: Ungerleider and Mishkin, “Two Cortical Visual Systems” (1982)

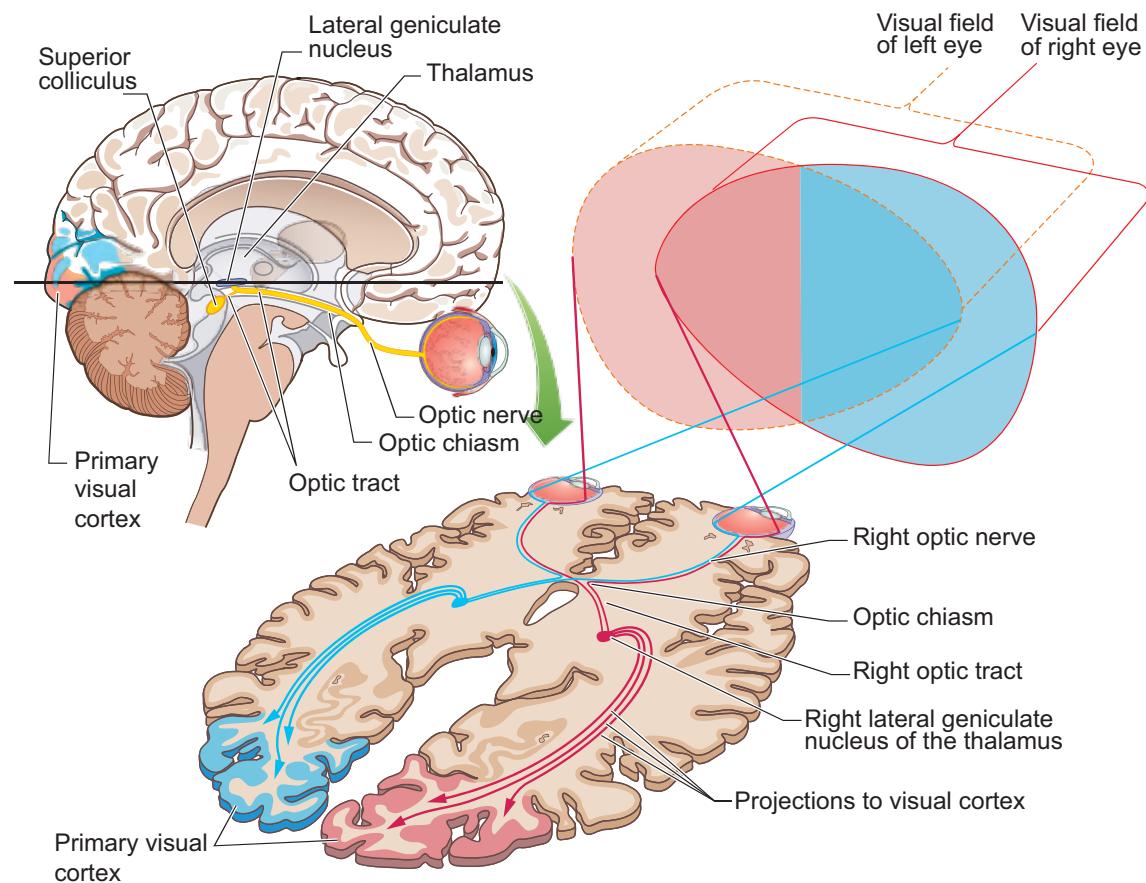
This section introduces the two visual systems hypothesis, first proposed by the neurologists Leslie Ungerleider and Mortimer Mishkin. The two visual systems hypothesis is



**Figure 3.3** The division of the left cerebral hemisphere into lobes.

important both because of the tools that were used to arrive at it (including the study of brain-damaged patients and experiments on monkeys) and because it illustrates a bottom-up, as opposed to top-down, way of studying the mind.

Ungerleider and Mishkin suggested that visual information does not take a single route from the primary visual cortex. Instead, the route the information takes depends upon the type of information it is. Information relevant to recognizing and identifying objects



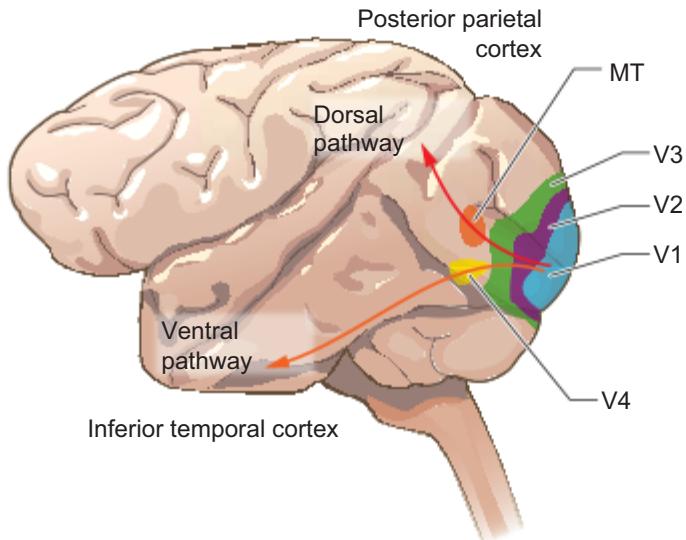
**Figure 3.4** The primary visual pathway. Note the contralateral organization, with information from the right side of space processed by the left side of the brain.

follows a *ventral route* (see Box 3.2) from the primary visual cortex to the temporal lobe, while information relevant to locating objects in space follows a *dorsal route* from the primary visual cortex to the posterior parietal lobe. The two routes are illustrated in Figure 3.5.

## BOX 3.2 Brain Vocabulary

Neuroscientists and neuroanatomists use an unfamiliar vocabulary for talking about the layout of the brain:

- Rostral** = at the front
- Caudal** = at the back
- Ventral** = at the bottom
- Dorsal** = at the top
- Ipsilateral** = same side
- Contralateral** = opposite side



**Figure 3.5** Image showing ventral (purple) and dorsal (green) pathways in the human visual system.

Ungerleider and Mishkin came to this conclusion after studying cognitive impairments due to brain damage and performing neuroanatomical experiments on monkeys. The neuroanatomical experiments were their distinctive contribution. By the time they were writing, there was already considerable evidence from brain-damaged patients that damage to the temporal and parietal lobes produced very different types of cognitive problem. Damage to the temporal cortex is associated with problems in identifying and recognizing objects, while damage to the parietal cortex tends to result in problems locating objects.

Evidence of this type has always been very important in working out the function of the different lobes (see Box 3.1 for a standard “division of labor” between the lobes). But being able to localize specific functions in this way falls a long way short of telling us the full story about the path that information takes in the brain. For that Ungerleider and Mishkin turned to experiments on monkeys.

The particular type of experiments that they carried out are called *cross-lesion disconnection experiments*. This is a methodology explicitly designed to trace the connections between cortical areas and so to uncover the pathways along which information flows. It addresses a fundamental problem with making inferences about the function and specialization of particular brain areas from what happens when those areas are damaged. Simply finding specific cognitive problems associated with damage to a specific brain region gives us no way of telling whether the impaired cognitive abilities are normally carried out by the damaged brain region itself, or by some other brain region that crucially depends upon input from the damaged brain region. Solving this problem cannot be done simply by observing the results of brain damage. Precise surgical intervention is required, in the form of targeted removal of specific brain areas to uncover the connections between them.

The cross-lesion disconnection experiments exploit the fact that the cerebrum is divided into two hemispheres, with duplication of the principal cortical areas. Suppose that investigators think that they have identified a cortical pathway that connects two cortical areas.

They can remove the area assumed to be earlier in the pathway from one hemisphere and the area assumed to be later from the other hemisphere.

Ungerleider and Mishkin, for example, hypothesized that there is a pathway connecting the primary visual cortex and the inferior temporal area, and so they performed surgery in monkeys to remove the primary visual cortex from one hemisphere in monkeys and the inferior temporal area from the other hemisphere. This destroyed the postulated pathway in each hemisphere. However, because the hemispheres can communicate through a large bundle of fibers known as the *corpus callosum* (illustrated in Figure 3.2), it turned out that there was little or no loss of function in the monkeys.

So, for example, it is well documented that monkeys who have had their inferior temporal cortex removed from both hemispheres are severely impaired on basic pattern discrimination tasks. But these pattern discrimination tasks were successfully performed by monkeys with primary visual cortex removed from one hemisphere and inferior temporal cortex from the other. Cutting the corpus callosum, however, reduced performance on those pattern discrimination tasks to chance and the monkeys were unable to relearn it.

Using experiments such as these (in addition to other types of neurophysiological evidence), Ungerleider and Mishkin conjectured that information relevant to object identification and recognition flows from the primary visual cortex to the inferior temporal cortex via areas in the occipital lobe collectively known as the *prestriate cortex*. They called this the ventral pathway.

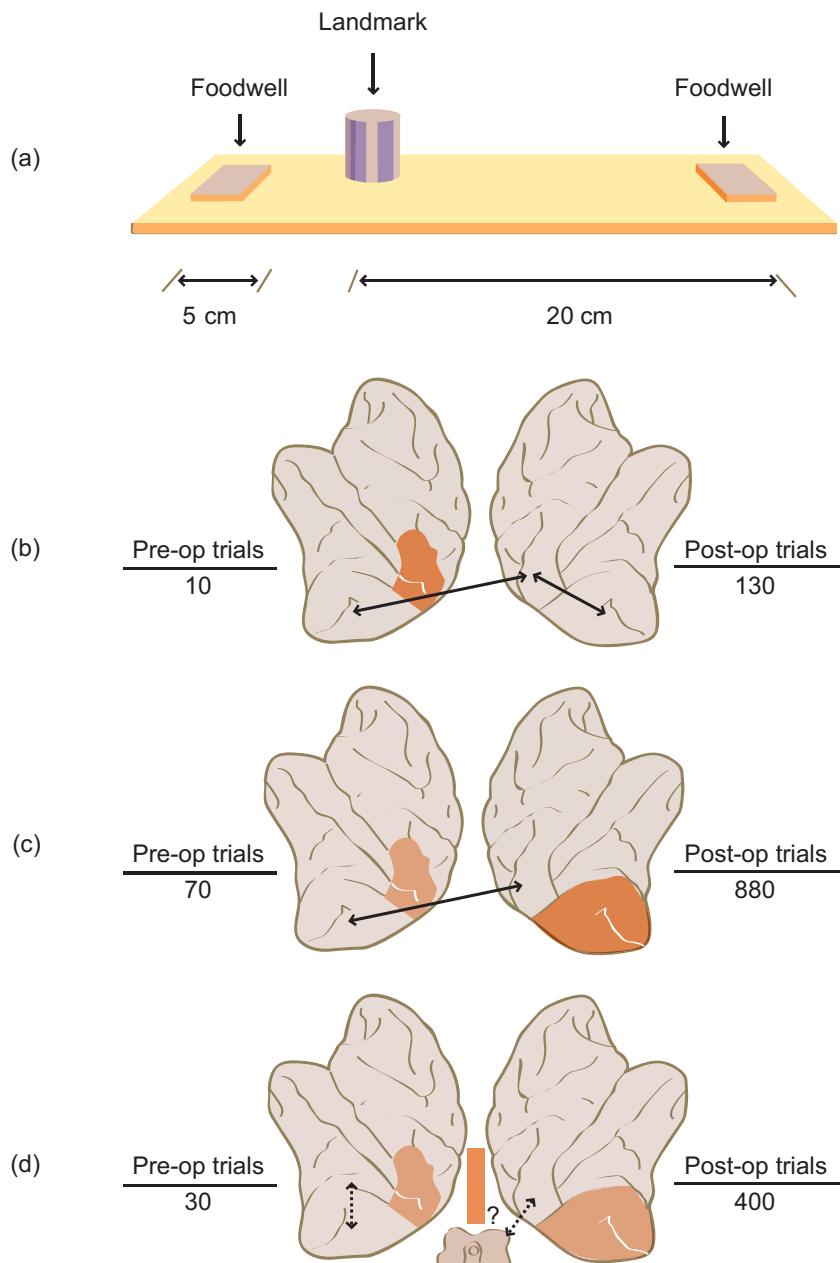
Ungerleider and Mishkin identified a completely different pathway (the *dorsal pathway*) leading from the primary visual cortex to the posterior parietal lobe. Once again they used cross-lesion disconnection experiments. In this case the task was the so-called landmark task, illustrated in the top left part of Figure 3.6.

In the landmark task monkeys are trained to choose food from one of two covered foodwells, depending on its proximity to a striped cylinder. The striped cylinder is moved at random and the task tests the monkey's ability to represent the spatial relation between the striped cylinder and the two foodwells.

The basic methodology was the same as for the experiments on the visual recognition pathway. The surgery proceeded in three stages. In the first stage (b in Figure 3.6) the posterior parietal cortex was removed from one side. The second stage (c) removed the primary visual cortex on the opposite side. The final stage (d) was a *transection* (severing) of the corpus callosum.

As indicated in Figure 3.6, the monkeys were tested on the landmark task both before and after each stage. However, the impairments on the landmark task were much more complicated than in the earlier experiments. The numbers in Figure 3.6 indicate the number of trials required to train the monkeys to a 90 percent success rate on the landmark task. So, for example, prior to the first stage of the surgery the average number of training trials required was ten. After lesion of the posterior parietal cortex the number of training trials went up to 130.

One interesting feature of these experiments is that the most severe impairment was caused by the second stage in the surgery, the removal of the primary visual cortex (in contrast to the other experiments on the visual recognition pathway, where severe impairments appeared only with the cutting of the corpus callosum). Ungerleider and Mishkin concluded from this that the posterior parietal cortex in a given hemisphere does not



**Figure 3.6** Design and results of Ungerleider and Mishkin's cross-lesion disconnection studies. (a) Landmark task. Monkeys were rewarded for choosing the covered foodwell located closer to a striped cylinder (the "landmark"), which was positioned on the left or the right randomly from trial to trial, but always 5 cm from one foodwell and 20 cm from the other. Training was given for thirty trials per day to a criterion of ninety correct responses in 100 consecutive trials. (b) Discrimination retention before and after first-stage lesion (unilateral posterior parietal;  $V = 3$ ); 10 preoperative trials and 130 postoperative trials. (c) Discrimination retention before and after second-stage lesion (contralateral striate;  $y = 3$ ); 70 preoperative and 880 postoperative trials. (d) Discrimination retention before and after third-stage lesion (corpus callosum;  $N = 3$ ); 30 preoperative and 400 postoperative trials. At each stage, the lesion is shown in dark brown and the lesions of prior stages in light brown. Arrows denote hypothetical connections left intact by lesions. (Adapted from Ungerleider and Mishkin 1982)

depend much upon information about the ipsilateral visual field (see Box 3.2) from the opposite hemisphere's primary visual cortex.

This raises the following intriguing possibility, since it is known that each hemisphere is specialized for the contralateral region of space. It may be that the posterior parietal cortex in each hemisphere is specialized for processing information about the opposite region of space. This would mean, for example, that the left posterior parietal cortex processes information about the layout of space on the perceiver's right-hand side.

This could be particularly important for thinking about the neurological disorder of *unilateral spatial neglect*. Patients with this disorder typically "neglect" one-half of the space around them, eating food from only one side of the plate and describing themselves as unaware of stimuli in the neglected half of space. Unilateral spatial neglect typically follows damage to the posterior parietal cortex in one hemisphere (most often the right) and the neglected region is contralateral to the damage (so that, most often, the left-hand side of space is neglected).

The visual systems hypothesis was a very important step in mapping out the *connectivity* of the brain. Ungerleider and Mishkin's basic distinction between the "what" system (served by the ventral pathway) and the "where" system (served by the dorsal pathway) has been refined and modified by many researchers (see the references in the Further Reading section of this chapter). However, the idea that there is no single pathway specialized for processing visual information, but instead that visual information takes different processing routes depending upon what type of information it is, has proved very enduring. From the perspective of cognitive science, the significance of the two visual systems hypothesis is that it exemplifies in a particularly clear way the bottom-up study of how information is processed in the mind.

There are recognizable affinities between what Ungerleider and Mishkin were doing, on the one hand, and the top-down approach of cognitive scientists such as Marr, on the other. So, for example, both are concerned with identifying distinct processing systems in terms of the functions that they perform. The real difference comes, however, with how they arrive at their functional analyses.

For Marr, the primary driver is top-down thinking about the role of visual processing within the overall organization of cognition and the behavior of the organism. For Ungerleider and Mishkin, the primary driver is thinking that starts at what Marr would term the implementational level. Instead of abstracting away from details of the channels and pathways between neural systems along which information processing flows, Ungerleider and Mishkin started with those channels and pathways and worked upward to identifying distinct cognitive systems carrying out distinct cognitive functions.



**Exercise 3.2** Make as detailed a list as you can of similarities and differences between the top-down and bottom-up approaches to studying the organization of the mind.

### 3.3 Extending Computational Modeling to the Brain

The historical antecedents of neutrally inspired computational models go back to the 1940s, and in particular to the work of Warren McCulloch and Walter Pitts (discussed in more detail in Chapter 5). But an important modern landmark was the publication in



1986 of a very influential two-volume collection of papers by Rumelhart, McClelland, and the PDP Research Group.

## A New Set of Algorithms: Rumelhart, McClelland, and the PDP Research Group, *Parallel Distributed Processing: Explorations in the Microstructure of Cognition* (1986)

The papers in the collection proposed and pursued a new set of abstract mathematical tools for modeling cognitive processes. These models are sometimes called *connectionist* networks and sometimes *artificial neural networks* (terms that we will use interchangeably). They abstract away from many biological details of neural functioning. But they manage to capture some of the key features of how the brain works. We will be looking in much more detail at artificial neural networks in later chapters (particularly Chapter 5). Here we will simply give a brief sketch of some key features.

First, connectionist networks display *parallel processing*. An artificial neural network contains a large number of units (artificial neurons). Each unit has a varying level of activation, typically represented by a real number between -1 and 1. The units are organized into layers with the activation value of a given layer determined by the activation values of all the individual units. The simultaneous activation of these units, and the consequent spread of activation through the layers of the network, governs how information is processed within the network. The processing is parallel because the flow of information through the network is determined by what happens in all of the units in a given layer – but none of those units are connected to each other.

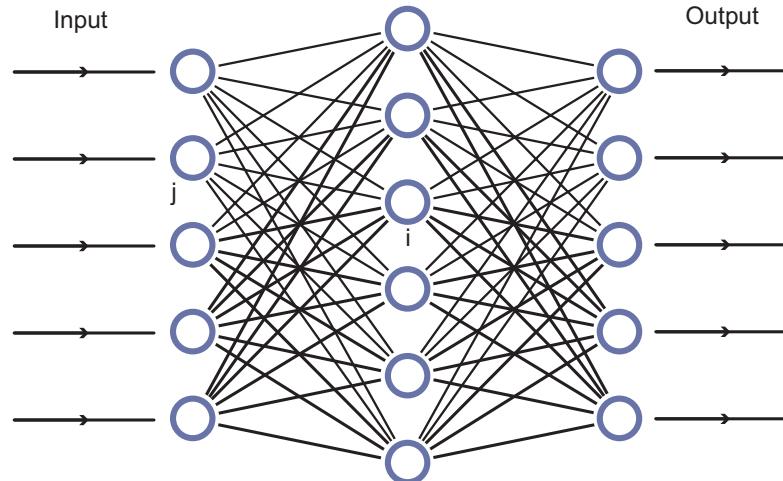
Second, each unit in a given layer has connections running to it from units in the previous layer (unless it is a unit in the input layer) and will have connections running forward to units in the next layer (unless it is a unit in the output layer). The pattern of connections running to and from a given unit is what identifies that unit within the network. The strength of the connections (the *weight* of the connection) between individual neurons varies. In fact, neural networks learn by modifying their weights.

Third, there are no intrinsic differences between one unit and another. The differences lie in the connections holding between that unit and other units.

Finally, most artificial neural networks are trained, rather than programmed. They are generally constructed with broad, general-purpose learning algorithms that work by changing the connection weights between units in a way that eventually yields the desired outputs for the appropriate inputs. These algorithms work by changing the weights of the connections between pairs of neurons in adjacent layers in order to reduce the “mistakes” that the network makes.

Figure 3.7 illustrates a generic neural network with three layers of units. The first layer is made up of input units, which receive inputs from sources outside the network. The third layer is made up of output units, which send signals outside the network. The middle layer is composed of what are called hidden units.

Hidden units only communicate with units within the network. They are the key to the computational power of artificial neural networks. Networks without hidden units can



**Figure 3.7** A generic three-layer connectionist network (also known as an artificial neural network). The network has one layer of hidden units. (Adapted from McLeod, Plunkett, and Rolls 1998)

only perform very limited computations. The illustrated network only has one layer of hidden units, but in fact networks can be constructed with as many layers as required. (More details coming up in Chapter 5.)

It takes a long time to train a network. Typically the network starts with randomly assigned weights. It is then given a training series of input patterns of activation, each of which is associated with a target output pattern of activation. The input patterns are presented. Differences between the actual output pattern and the target output pattern result in changes to the weights. (This is what the learning algorithm does – adjust the weights in order to reduce the difference between actual and desired output.)

This training process continues until errors have diminished almost to zero, resulting in a distinctive and stable pattern of weights across the network. The overall success of a network can be calculated by its ability to produce the correct response to inputs on which it has not been trained. The next subsection illustrates the sort of task that a network can be trained to do with a justly celebrated example.

### Pattern Recognition in Neural Networks: Gorman and Sejnowski, "Analysis of Hidden Units in a Layered Network Trained to Identify Sonar Targets" (1998)

Artificial neural networks are particularly suited for pattern recognition tasks. Here is a classic example. Consider the task of identifying whether a particular underwater sonar echo comes from a submerged mine, or from a rock. There are discriminable differences between the sonar echoes of mines and rocks, but there are equally discriminable differences between the sonar echoes from different parts of a single mine, or from different



parts of a single rock. It is no easy matter to identify reliably whether a sonar echo comes from a mine or from a rock. Human sonar operators can do so reasonably well (after a considerable amount of practice and training), but it turns out that artificial neural networks can perform significantly better than humans.

The first problem is coding the external stimulus as a pattern of activation values. The external stimuli are sonar echoes from similarly shaped and sized objects known to be either mines or rocks. In order to “transform” these sonar echoes into a representational format suitable for processing by the network, the sonar echoes are run through a spectral analyzer that registers their energy levels at a range of different frequencies. This process gives each sonar echo a unique “fingerprint” to serve as input to the network. Each input unit is dedicated to a different frequency and its activation level for a given sonar echo is a function of the level of energy in the relevant sonar echo at that frequency. This allows the vector of activation values defined over the input units to reflect the unique fingerprint of each sonar echo.

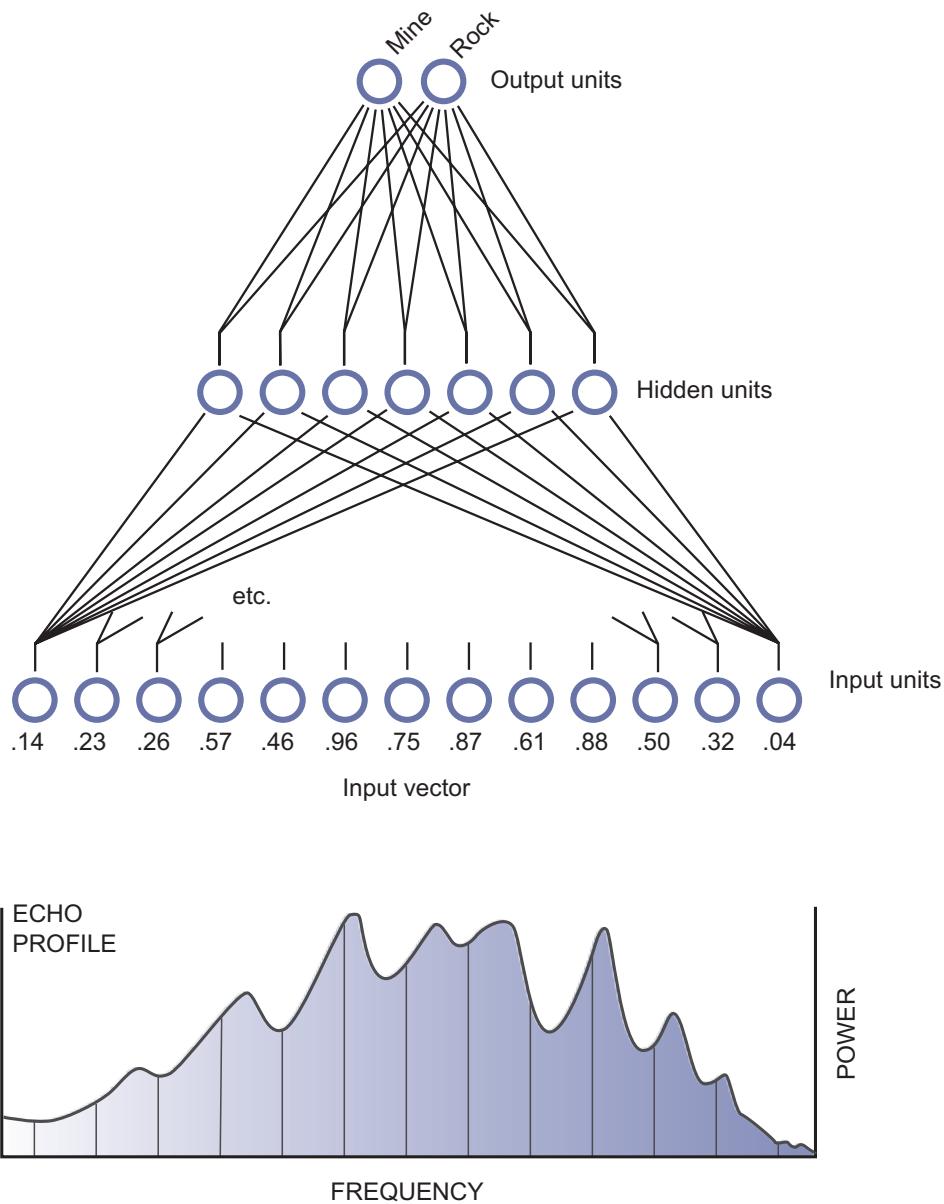
The neural network developed by Paul Gorman and Terrence Sejnowski to solve this problem contains sixty input units, corresponding to the sixty different frequencies at which energy sampling was carried out, and one layer of hidden units. Since the job of the unit is to classify inputs into two groups, the network contains two output units – in effect, a rock unit and a mine unit. The aim of the network is to deliver an output activation vector of  $<1,0>$  in response to the energy profile of a rock and  $<0,1>$  in response to the energy profile of a mine. Figure 3.8 is a diagrammatic representation of Gorman and Sejnowski’s mine/rock network.

The mine detector network is a standard feedforward network (which means that activation is only ever spread forward through the network) and is trained with the back-propagation learning algorithm (explained in Chapter 5). Although the network receives information during the training phase about the accuracy of its outputs, the only memory it has of what happened in early sessions is the particular patterns of weights holding across the network. Each time the network comes up with a wrong output (a pattern of  $<0.83, 0.17>$  rather than  $<1.0>$ , for example, in response to a rock profile), the error is propagated backward through the network and the weights adjusted to reduce the error. Eventually the error at the output units diminishes to a point where the network can generalize to new activation patterns with a 90 percent level of accuracy.

The mine/rock detection task is a paradigm of the sort of task for which neural networks are best known and most frequently designed. The essence of a neural network is pattern recognition. But many different types of cognitive ability count as forms of pattern recognition and the tools provided by artificial neural networks have been used to model a range of cognitive processes – as well as many phenomena that are not cognitive at all (such as predicting patterns in the movements of prices on the stock markets, valuing bonds, and forecasting demand for commodities).



**Exercise 3.4** Give examples of cognitive abilities that you think would lend themselves to being modeled by artificial neural networks.



**Figure 3.8** Gorman and Sejnowski's mine/rock detector network. (Adapted from Gorman and Sejnowski 1988)

## 3.4

### Mapping the Stages of Lexical Processing

In this section we turn back from thinking about computational modeling in the abstract to thinking about how the direct study of the brain can help cognitive scientists to formulate and decide between different models. Earlier in this chapter we looked at how neurological experiments on monkeys have been used to identify the channels and pathways along which visual information flows. We turn now to a different set of techniques that have become an increasingly important part of the cognitive scientist's tool kit.



## Functional Neuroimaging with PET: Petersen, Fox, Posner, and Mintun, "Positron Emission Tomographic Studies of the Cortical Anatomy of Single-Word Processing" (1988)

Functional neuroimaging allows brain activity to be studied noninvasively. No surgery is required and subjects can be studied while they are actually performing experimental tasks.

There are different types of functional neuroimaging. The first experiments we will look at use the technique known as *positron emission tomography* (better known under its acronym PET). We will be looking at fMRI (*functional magnetic resonance imaging*) in the next section.

The basic idea behind the PET technology (as with functional neuroimaging in general) is to study the function of different brain areas by measuring blood flow in the brain. We can work out which brain areas are involved in carrying out particular cognitive tasks by identifying the areas to which blood is flowing. The distinctiveness of PET is that it provides a safe and precise way of measuring short-term blood flow in the brain. Subjects are given (typically by injection) a small quantity of water containing the positron-emitting radioactive isotope oxygen-15 ( $^{15}\text{O}$ ). The radioactive water accumulates in the brain in direct proportion to the local blood flow, so that areas to which the most blood is flowing will show the greatest concentration of  $^{15}\text{O}$ .

The PET scanner tracks the progress of the radioactive water through the brain (for about a minute, before the radioactive isotope decays to a nonradioactive atom). This provides an indirect, but highly reliable, measure of blood flow in the brain, and hence a way of telling which brain regions are active during the minute after administering the water. If subjects are carrying out particular experimental tasks during that time, then the PET technology gives scientists a tool for identifying which brain regions are actively involved in carrying out that task.

However, simply identifying which brain regions have blood flowing to them while a particular task is being performed is not enough to tell us which brain regions are actively involved in carrying out the task. There may be all sorts of activity going on in the brain that are not specific to the particular experiment that the subject is performing. The art in designing PET experiments is finding ways to filter out potentially irrelevant, background activity. The experiments we will be focusing on, carried out by Steve Petersen and collaborators at Washington University in St. Louis, provide a very nice illustration of how this sort of filtering can be done – and of how careful experimental work can refine information-processing models.

## Petersen, Fox, Posner, and Mintun, "Positron Emission Tomographic Studies of the Cortical Anatomy of Single-Word Processing" (1988)

Petersen and his colleagues were studying how linguistic information is processed in the human brain. They started with individual words – the basic building blocks of language.

Many different types of information are relevant to the normal course of reading, writing, or conversing. There is visual information about the shape and layout of the word, as well as auditory information about how the word sounds and *semantic* information about what the word means. The interesting question is how these different types of information are connected together. Does silently reading a word to oneself involve processing information about how the word sounds? Does simply repeating a word involve recruiting information about what the word means?

The two leading information-processing models of single-word processing (often called *lexical access*) answer these two questions very differently. Within neurology the dominant model, derived primarily from observing brain-damaged patients, holds that the processing of individual words in normal subjects follows a single, largely invariant path. The information-processing channel begins in the sensory areas. Auditory information about how the word sounds is processed in a separate brain region from information about the word's visual appearance. According to the neurological model, however, visual information about the word's appearance needs to be phonologically recoded before it can undergo further processing. So, in order to access semantic information about what a written word means, the neurological model holds that the brain needs to work out what the word sounds like. Moreover, on this model, semantic processing is an essential preliminary to producing phonological motor output. So, for example, reading a word and then pronouncing it aloud involves recruiting information about what the word means.



### **Exercise 3.5** Draw a flowchart illustrating the distinct information-processing stages in single-word processing according to the neurological model.

The principal alternative to the neurological model is the cognitive model (derived primarily from experiments on normal subjects, rather than from studies of brain-damaged patients). The neurological model is *serial*. It holds that information travels through a fixed series of information-processing "stations" in a fixed order. In contrast, the cognitive model holds that lexical information processing is *parallel*. The brain can carry out different types of lexical information processing at once, with several channels that can feed into semantic processing. Likewise, there is no single route into phonological output processing.

Petersen and colleagues designed a complex experiment to determine which model reflects more accurately the channels of lexical information processing in the brain. The basic idea was to organize the experimental conditions hierarchically, so that each condition could tap into a more advanced level of information processing than its predecessor. Each level involved a new type of information-processing task. Successfully carrying out the new task required successfully carrying out the other tasks lower in the hierarchy. What this means is that by looking at which *new* brain areas are activated in each task we can identify the brain areas that are specifically involved in performing that task – and we can also see which brain areas are *not* involved.



The baseline condition was simply asking subjects to focus on a fixation point (a small cross-hair) in the middle of a television screen. The point of asking the subjects to do this was to identify what is going on in the brain when subjects are visually attending to something that is not a word. The second condition measured brain activity while subjects were passively presented with words flashed on the screen at a rate of forty words per minute. The subjects were not asked to make any response to the words. In a separate condition the same words were spoken to the subjects.

Combining the results from these two different conditions allowed Petersen and his colleagues to work out which brain areas are involved in visual and auditory word perception. The key to doing this is to subtract the image gained from the first condition from the image derived from the second condition. The image of brain activity while fixating on the cross-hair acts as a control state. In principle (and we will look much more closely at some of the methodological difficulties in functional neuroimaging in Chapter 9), this allows us to filter out all the brain activation that is responsible for sensory processing in general, rather than word perception in particular.

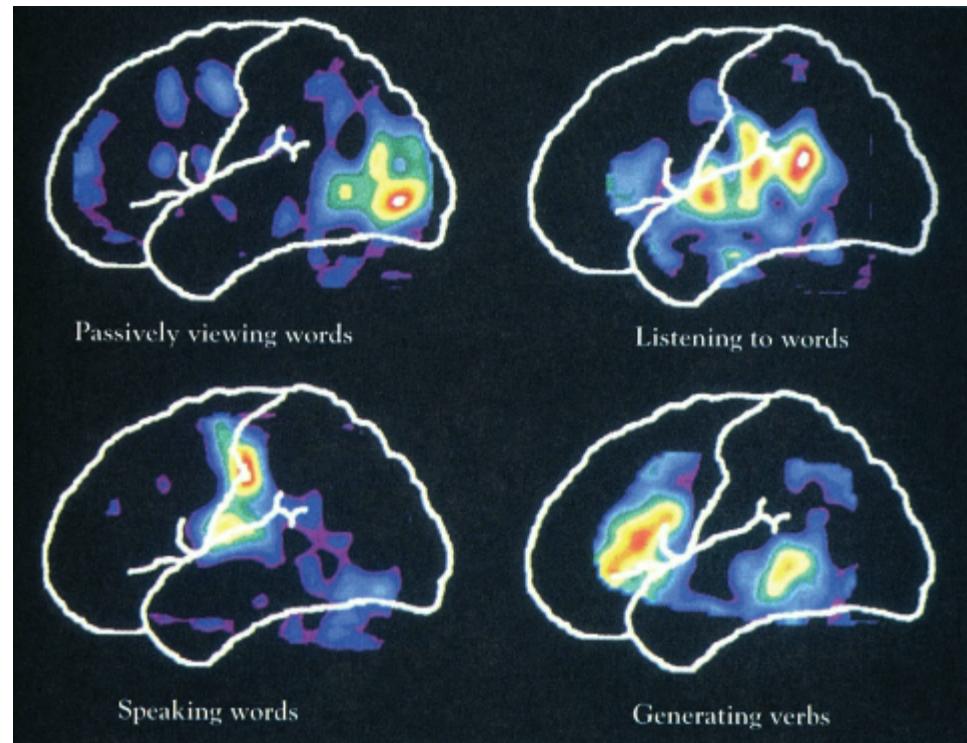
The third and fourth levels of the experimental hierarchy measured brain activation during more complex tasks. The aim here was to trace the connections between initial sensory processing and the semantic and output processing that takes place further “downstream.” In the third condition subjects were asked to say out loud the word appearing on the screen. Subtracting the resulting image from the word perception image allowed Petersen and his colleagues to calculate which brain areas are involved in speech production. Finally, the highest level of the experimental hierarchy involved a task that clearly requires semantic processing. Here the subjects were presented with nouns on the television monitor and asked to utter an associated verb. So, for example, a subject might say “turn” when presented with the word “handlebars.” As before, Petersen and his colleagues argued that subtracting the image of brain activation during this semantic association task from the image obtained from the speech production task would identify the brain areas involved in semantic processing.



### **Exercise 3.6** Make a table to show the different levels in the hierarchy and the aspects of single-word processing that they are intended to track.

Statistical comparison of the brain images in the different stages of the experiment produced a number of striking results. As we see in Figure 3.9, each of the tasks activated very different sets of brain areas. (The areas with the maximum blood flow are colored white, followed in decreasing order by shades of red, yellow, green, blue, and purple.)

Moreover, the patterns of activation seemed to provide clear evidence against the neurological model. In particular, when subjects were asked to repeat visually presented words, there was no activation of the regions associated with auditory processing. This suggested to Petersen and his colleagues that there is a direct information pathway from the areas in the visual cortex associated with visual word processing to the distributed network of areas responsible for articulatory coding and motor



**Figure 3.9** Images showing the different areas of activation (as measured by blood flow) during the four different stages in Petersen et al.'s (1988) lexical access studies. (From Posner and Raichle 1994)

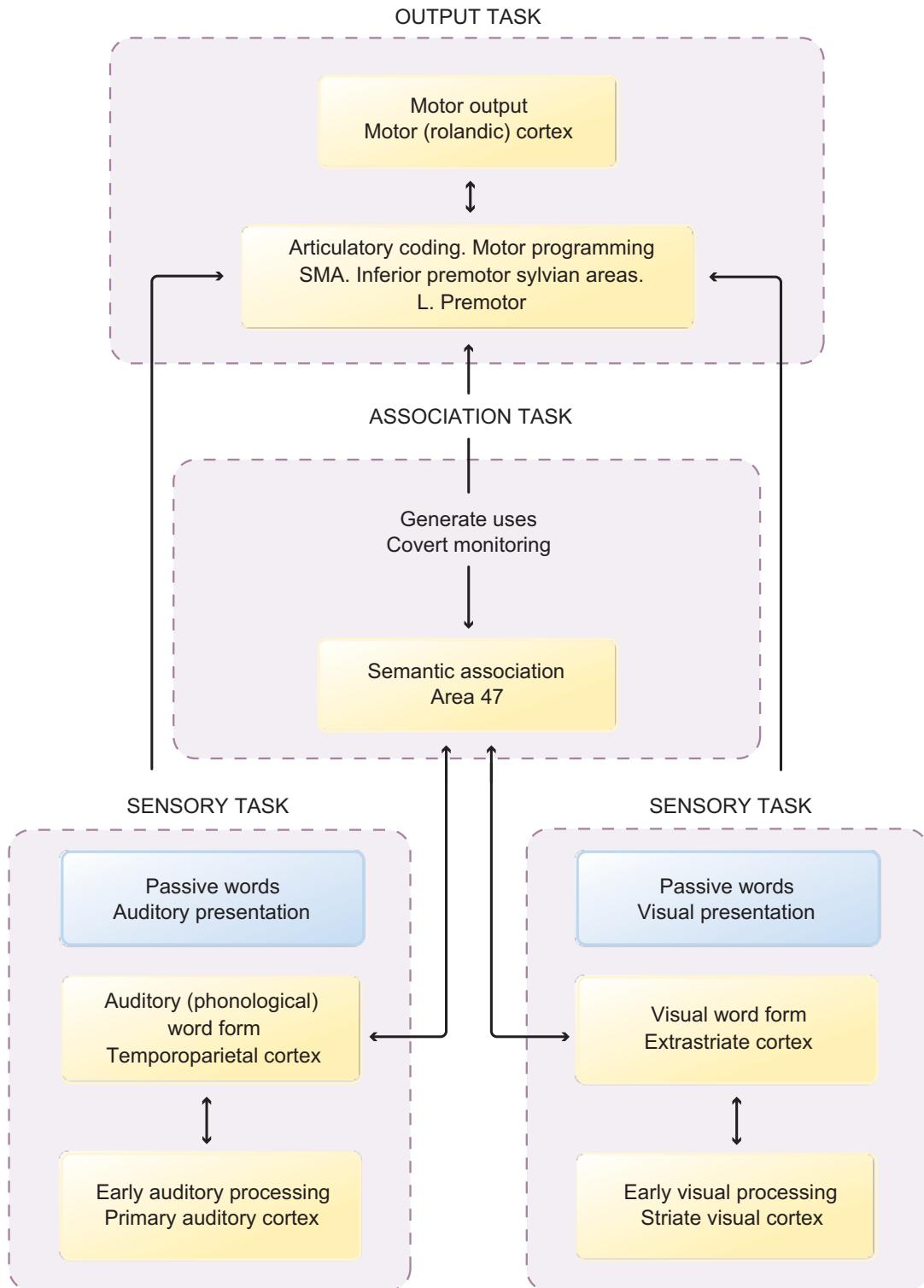
programming, coupled with a parallel and equally direct pathway from the areas associated with auditory word processing. Moreover, the areas associated with semantic processing (those identified in the condition at the top of the hierarchy) were not involved in any of the other tasks, suggesting that those direct pathways did not proceed via the semantic areas.

The situation can most easily be appreciated in an information-processing diagram. Figure 3.10 is drawn from a paper by Petersen and collaborators published in the journal *Nature* in 1988. Unlike many information-processing flowcharts, this one is distinctive in that it identifies the particular brain areas that are thought to carry out each distinct stage. This is not an accident. It reflects how the information-processing model was reached – on the basis of direct study of the brain through PET scan technology.

## 3.5

### Studying Memory for Visual Events

There are two principal technologies in functional neuroimaging. In Section 3.4 we looked at the PET technology, which measures cerebral blood flow by tracking the movement of radioactive water in the brain. A newer, and by now dominant, technology is functional magnetic resonance imaging (fMRI).



**Figure 3.10** A flowchart relating areas of activation in Petersen et al.'s 1988 study to different levels of lexical processing. The dashed boxes outline the different subtraction. The solid boxes outline possible levels of coding and associated anatomical areas of activation. (From Petersen et al. 1988)



## Functional Neuroimaging with fMRI

The standard background assumption in fMRI (as with PET) is that blood flow to a particular region of the brain increases when cellular activity in that region increases. This increase in blood flow produces an increase in oxygen. Since the supply of oxygen is greater than the demand, the blood oxygen level increases in a brain region that is undergoing increased cellular activity. An fMRI scanner creates a powerful magnetic field, which can detect increases in blood oxygen, since oxygenated and deoxygenated blood have different properties.

The difference between oxygenated and deoxygenated blood is known as the BOLD (blood oxygen level dependent) contrast. Functional magnetic resonance imaging measures the BOLD signal. (As we will see in the next section, there are different ways of thinking about what the BOLD signal tells us about the brain.)

Early fMRI experiments used a similar experimental design to those used in PET studies, such as the Peterson et al. experiments described in Section 3.4. This is known as a blocked design. In a blocked design, experiment subjects perform a task for an extended period of time. The task might be focusing on a fixation point, or repeated visually presented words. Using extended blocks maximizes the changes in blood oxygen level and so gives a stronger signal. This in turn makes it easier to compare and contrast signals in different conditions using subtraction methods.

One reason for the popularity of blocked designs in early fMRI experiments is that the BOLD hemodynamic response is delayed in its onset and takes a while to develop. So, for example, for a neural event lasting one second, it takes two seconds for the hemodynamic response to start to develop, and the development takes 10-12 seconds. So, early researchers concluded that these aspects of the BOLD signal required fMRI experiments to use the same kind of experimental design as PET experiments.

In the early 1990s, however, a new type of experimental design emerged for fMRI experiments. This is usually termed *event-related fMRI*. It is in some respect similar to the experimental design used in many EEG (electroencephalography) experiments, which we will look at in more detail in Chapter 9. The basic idea of event-related fMRI is to measure the BOLD signal associated with individual rapid occurring neural events, even though these events elicit overlapping hemodynamic responses. Event-related fMRI is possible because the hemodynamic response measured through fMRI behaves (to a first approximation) like a linear system.

What this means is that the hemodynamic response for a given event in a series basically adds on proportionally to the hemodynamic response for earlier events in the series. So, if you can measure the change in the BOLD signal associated with six task-events, each two seconds in length, the result will be close to the BOLD signal associated with a single task-event, lasting for 12 seconds. Because of this linearity, even though the hemodynamic responses for those six events are all overlapping, standard statistical techniques can be used to work backward from the overall change in the BOLD signal for the series to the particular change in the BOLD signal for each event. And it is also



possible, although more complicated, to do this when the events in the series are different from each other.

Using event-related fMRI, neuroscientists are able to study the BOLD signal produced by short duration events, and also to disentangle the separate components of complex tasks. This is a very significant step beyond the blocked-design paradigm (although that remains useful for various areas of neuroimaging, such as the neuroimaging for language). To illustrate the power of event-related fMRI, we will look at an important early study on memory for visual events.

## **Brewer, Zhao, Desmond, Glover, and Gabrieli, "Making Memories: Brain Activity That Predicts How Well Visual Experience Will Be Remembered" (1998)**

This important paper was one of the first fMRI studies to use an event-related experimental design. It was trying to get at something that could not be uncovered using a blocked design, because it depended on information about changes in the BOLD signal brought about by very specific and short-lived neural events. Brewer and his colleagues were interested in exploring whether there are any neural markers predicting how well specific visual experiences would be remembered. Are there any areas in the brain where activity in those areas would predict whether they would be remembered well, less well, or forgotten?

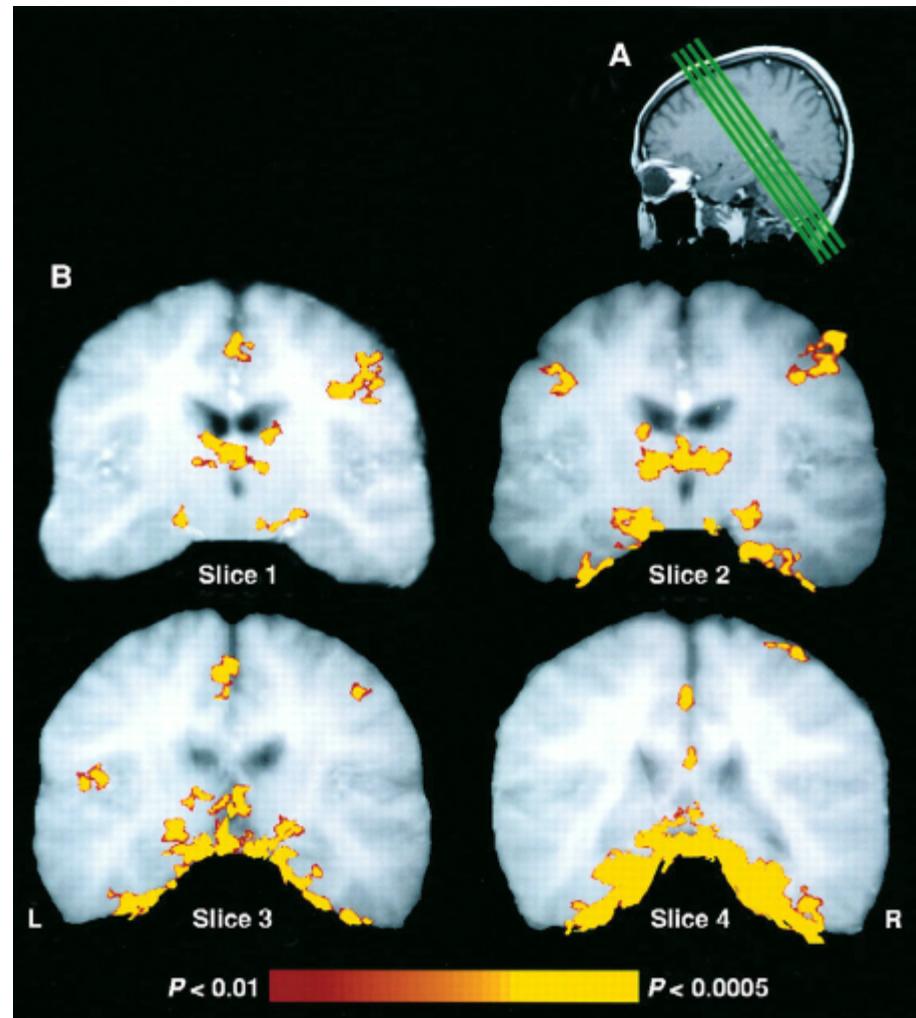
It is important to realize what this question is *not* asking. It is not asking which areas in the brain are involved in memory. That was already fairly well known from studies of brain-damaged patients. So, for example, there is considerable evidence that damage to the medial temporal lobe causes global amnesia, and that damage to specific parts of the frontal lobes can bring about different types of local amnesia. Damage to the left front lobes affect verbal memory, for example.

But looking at brain damage cannot distinguish between the different aspects of memory. It cannot tell us whether what is impaired is how experiences are encoded for memory; how memories are stored; or how they are retrieved. And it certainly cannot tell us anything about individual experiences, and how likely they are to be remembered. For that you need to be able to identify the specific hemodynamic response generated by specific experiences – in other words, you need an event-related design.



### **Exercise 3.7 Explain in your own words the difference between a blocked design and an event-related design.**

To get at this question about how individual experiences are encoded in memory, the experimenters showed subjects in an fMRI scanner ninety-six color pictures of indoor and outdoor scenes over four trials. The pictures were selected to be broadly comparable in complexity and visual quality, so that they would all make roughly the same general processing demands. All that the subjects were asked to do while in the scanner was to identify for each picture whether it was an indoor or outdoor scene. Then, 30 minutes later,

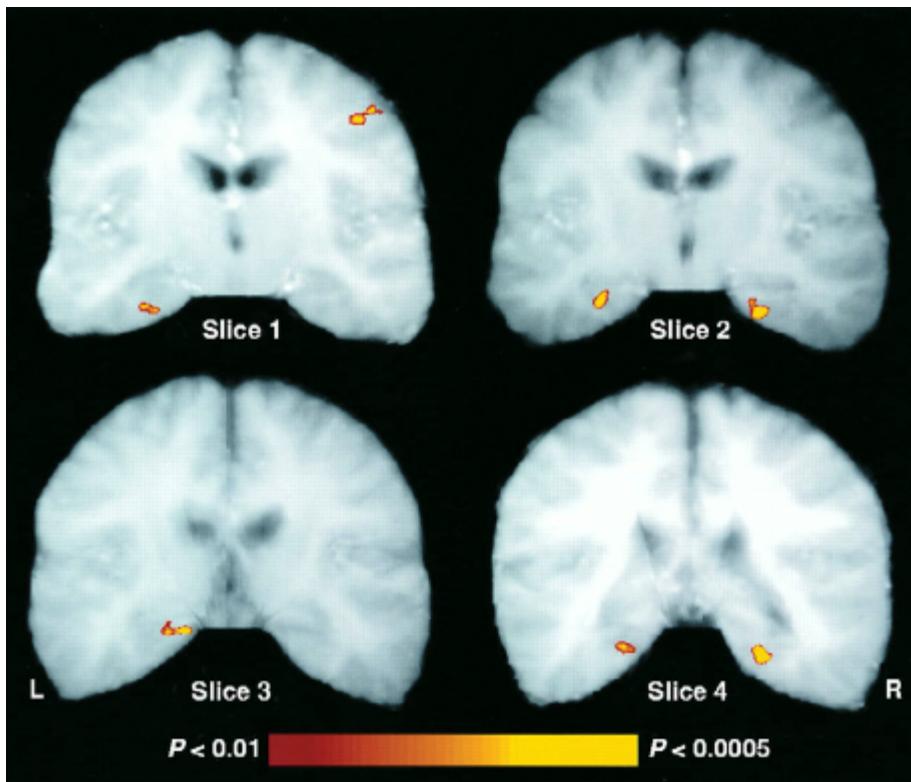


**Figure 3.11** Neural area showing activity when subjects looked at pictures.

the subjects were given an unanticipated memory test, presented with 128 pictures, including the 96 shown to them in the scanner, and asked to identify which they had seen before, and to identify how confident they were.

The memory test allowed the experimenters to classify how well each subject remembered each of the originally displayed pictures – as well remembered, as familiar, or as forgotten. The next step was to try to find patterns of neural activity correlated with each of those three levels of memory performance. To that end, the experimenters constructed two maps of event-related activity. The first map showed brain areas where the hemodynamic response increased when subjects were looking at pictures (relative to activation levels during fixation). So, this map measured the BOLD response generated by visual experiences. The first map is illustrated in Figure 3.11.

As Figure 3.11 shows, many areas are involved in visual experience. But which ones are responsible for encoding visual experiences into memory?



**Figure 3.12** Neural areas where activation is correlated with levels of memory performance.

The second map answered this question. It showed brain areas where activation levels were correlated with the levels of memory performance for individual events. So, for example, a brain area would appear on the second map if it showed high levels of increased activation during visual experiences of pictures that were subsequently well remembered, medium levels of increased activation during visual experiences that were subsequently judged familiar, and low levels of increased activation during visual experiences that were not remembered.

The second map is depicted in Figure 3.12.

Comparing Figures 3.11 and 3.12 shows that encoding visual experiences into memory is highly localized. Of the many areas implicated in visual experience, only two predict how well those visual experiences will be remembered. The first area is the parahippocampal cortex (in both hemispheres). The parahippocampal cortex is part of the medial temporal lobe. The second area is the dorsolateral prefrontal cortex, only in the right hemisphere.

Quite apart from the intrinsic interest of these results for the study of memory, they are an excellent illustration of the power of event-related designs for fMRI. The event-related design makes it possible to identify the hemodynamic response generated by each individual visual experience, and without that it would be impossible to identify the areas that predict which experiences will be remembered – as opposed to identifying the areas generally responsible for visual experience and/or memory.

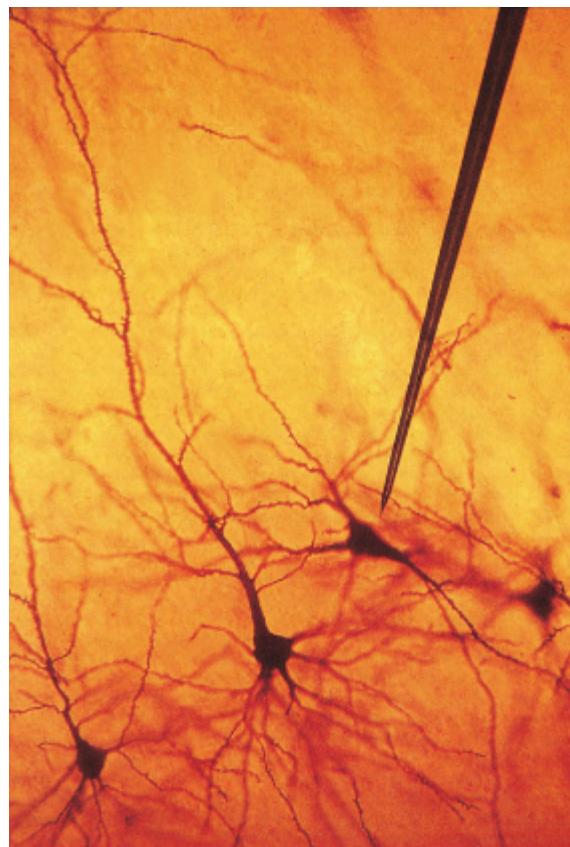
 3.6

## The Neural Correlates of the BOLD Signal

When one is looking at brightly colored pictures communicating the results of PET or fMRI scans it is only too easy to forget that relatively little is known about the relation between what those scans measure and the cognitive activity that is going on while the measurements are being made. It is only in the very recent past that progress has been made on building a bridge between functional neuroimaging and neurophysiology.

As we saw in Section 3.5, fMRI measures the BOLD contrast. But what does the BOLD contrast measure? In some sense the BOLD contrast has to be an index of cognitive activity—since it is known that cognitive activity involves increased activity in populations of neurons, which in turn results in increased oxygen levels and hence in a more pronounced BOLD contrast. But what sort of neuronal activity is it that generates the BOLD contrast?

Neuroscientists study the behavior of individual neurons through single-cell recordings (to be discussed in more detail in Chapter 9). Microelectrodes can be inserted into the brains of animals (and also of humans undergoing surgery) and then used to record activity in individual cells while the animal performs various behavioral tasks. Figure 3.13 illustrates a microelectrode recording in the vicinity of a single neuron.



**Figure 3.13** A microelectrode making an extracellular recording. (Scientific American Library [W. H. Freeman 1995])



Experimenters can track the relation between the firing rates of individual neurons and where the animal's attention is directed. These are usually low-level properties, such as the reflectance properties of surfaces. But in some cases neurons seem to be sensitive to higher-level properties, firing in response to particular types of object and/or situations. The basic assumption is that individual neurons are "tuned" to particular environmental properties.

Since the salient property of individual neurons is their firing (or *spiking*) behavior, it is a natural assumption that the neural activity correlated with the BOLD contrast is a function of the firing rates of populations of neurons. In fact, this is exactly what was suggested by Geraint Rees, Karl Friston, and Christoph Koch in a paper published in 2000. They proposed that there is a linear relationship between the average neuronal firing rate and the strength of the BOLD signal – two variables are linearly related when they increase in direct proportion to each other, so that if one were to plot their relation on a graph it would be a straight line.

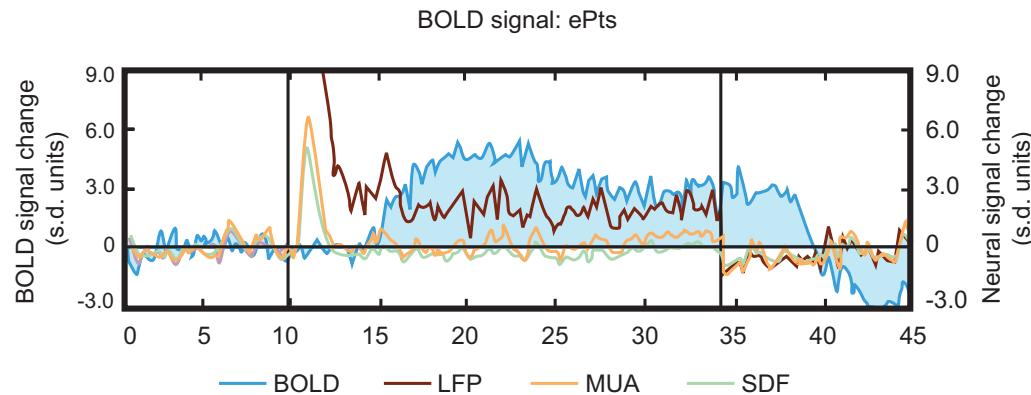
This conclusion was based on comparing human fMRI data with single-cell recordings from monkeys. In fact, their study seemed to show that each percentage increase in the BOLD contrast is correlated with an average per second increase of nine spikes per unit. If the Rees–Friston–Koch hypothesis is correct, then the BOLD response directly reflects the average firing rate of neurons in the relevant brain area, so that an increase in the BOLD contrast is an index of higher neural firing activity.

Neurons do more than simply fire, however. We can think of a neuron's firing as its *output*. When a neuron fires it sends a signal to the other neurons to which it is connected. This signal is the result of processing internal to the neuron. This processing does not always result in the neuron's firing. Neurons are selective. They fire only when the level of internal activity reaches a particular threshold. This means that there can be plenty of activity in a neuron even when that neuron does not fire.

We might think of this as a function of the *input* to a neuron, rather than of its output. A natural question to ask, therefore, is how cognitively relevant this activity is. And, given that we are thinking about the relation between neural activity and the BOLD contrast, we have a very precise way of formulating this question. We can ask whether the BOLD signal is correlated with the input to neurons, or with their output (as Rees, Friston, and Koch had proposed). This is exactly the question explored in a very influential experiment by Nikos Logothetis and collaborators.

## Logothetis, "The Underpinnings of the BOLD Functional Magnetic Resonance Imaging Signal" (2001)

Logothetis compared the strength of the BOLD signal against different measures of neural activity in the monkey primary visual cortex (see Section 3.2 for a refresher on where the primary visual cortex is and what it does). The team measured neural activity in an anesthetized monkey when it was stimulated with a rotating checkerboard pattern while in a scanner. In addition to using fMRI to measure the BOLD contrast, researchers used microelectrodes to measure both input neural activity and output neural activity. This is



**Figure 3.14** Simultaneous microelectrode and fMRI recordings from a cortical site showing the neural response to a pulse stimulus of 24 seconds. Both single- and multiunit responses adapt a couple of seconds after stimulus onset, with LFP remaining the only signal correlated with the BOLD response. (Adapted from Bandettini and Ungerleider 2001)

particularly challenging from an engineering point of view. Since an fMRI scanner generates a powerful magnetic field, the microelectrodes needed to be nonmagnetic.

At the output level they measured the firing rates both of single neurons and of small populations of neurons near the electrode tip (“near” here means within 0.2 mm or so). In Figure 3.14 these are labeled SDF (spike density function) and MUA (multiunit activity).

The *local field potential* (LFP) is an electrophysiological signal believed to be correlated with the sum of inputs to neurons in a particular area. It is also measured through a microelectrode, but the signal is passed through a *low-pass* filter that smooths out the quick fluctuations in the signal that are due to neurons firing and leaves only the low-frequency signal that represents the inputs into the area to which the electrode is sensitive (an area a few millimeters across).

The striking conclusion reached by Logothetis and his team is that the BOLD contrast is more highly correlated with the LFP than with the firing activity of neurons (either at the single-unit or multiunit level). This is nicely illustrated in the graph in Figure 3.14. In many cases, the LFP will itself be correlated with the firing activity of neurons (which is why Logothetis’s results are perfectly compatible with the results reached by Rees, Friston, and Koch). But, if Logothetis’s data do indeed generalize, then they show that when spiking activity and LFP are *not* correlated, the LFP is the more relevant of the two to the BOLD contrast.



## Summary

This chapter has explored the “turn to the brain” that took place in cognitive science during the 1980s. This involved the development of experimental paradigms for studying the information pathways in the brain from the bottom up. These experimental paradigms included lesion studies



on monkeys, as well as neuroimaging of human brains. It also involved the development of computational modeling techniques based on an idealized model of how neurons work.

The first example we looked at was the two visual systems hypothesis developed by Mishkin and Ungerleider primarily on the basis of monkey experiments. We then looked at the emergence of parallel distributed processing models in the 1980s and reviewed a famous application by Gorman and Sejnowski to the problem of distinguishing rocks from mines in sonar recordings. The last three sections focused on neuroimaging. We began with positron emission tomography (PET) and reviewed an influential set of experiments on single-word processing from Petersen, Fox, Posner, and Mintun. For many applications, PET has been superseded by functional magnetic resonance imaging (fMRI), not least because it allows event-related experimental designs, as opposed to the blocked designs used in PET and early fMRI studies. We then looked at the power of the event-related design in the context of experiments by Brewer, Zhao, Desmond, Glover, and Gabrieli on predicting how well visual experiences would be remembered. Finally, we turned to the neural correlates of the BOLD contrast, which is what fMRI measures directly. The BOLD contrast is a hemodynamic response, a function of blood oxygen levels. But what do blood oxygen levels tell us about neural activity? An elegant set of experiments by Nikos Logothetis tackle this question.

## Checklist

### Ungerleider and Mishkin's Two Visual Systems Hypothesis

- (1) The cross-lesion disconnection paradigm, coupled with various other anatomical and neurological methods, was used to identify two different information-processing pathways for visual information.
- (2) Both pathways start from the primary visual cortex.
- (3) Information relevant to object identification and recognition travels along the ventral pathway, from the primary visual cortex to the inferior temporal cortex via the prestriate cortex.
- (4) Information relevant to locating objects flows from the primary visual cortex to the posterior parietal lobe.

### Information Processing in Artificial Neural Networks

- (1) These networks are designed to reflect certain high-level features of how the brain processes information, such as its parallel and distributed nature.
- (2) The neuron-like units in artificial neural networks are organized into layers, with no connections between units in a single layer.
- (3) The overall behavior of the network is determined by the weights attached to the connections between pairs of units in adjacent layers.
- (4) Networks "learn" by adjusting the weights in order to reduce error.
- (5) Artificial neural networks are particularly suited to pattern recognition tasks, such as discriminating sonar echoes caused by mines from those caused by rocks.

### Functional Neuroimaging: PET and the Example of Single-Word Processing

- (1) PET allows brain activity to be studied noninvasively by measuring blood flow in the brain while subjects are performing particular cognitive tasks.

- (2) The paired-subtraction paradigm focuses on the brain activity specific to the task by subtracting out the activity generated by carefully chosen control tasks.
- (3) In studies of how single words are processed, experimenters constructed a four-level hierarchy of tasks of increasing complexity.
- (4) The patterns of activation they identified across the different tasks supported a parallel rather than a serial model of single-word processing.

### Functional Neuroimaging: Event-Related fMRI and Predicting How Well Visual Experiences Will Be Remembered

- (1) Functional magnetic resonance imaging (fMRI) measures levels of blood oxygen in the brain (the BOLD contrast/signal), and has superseded PET as a neuroimaging tool for many applications.
- (2) Event-related fMRI allows researchers to study the BOLD signal associated with individual neural events, unlike the blocked-design standard used in PET imaging.
- (3) Event-related fMRI works because, even though the BOLD response is delayed and takes time to develop, it behaves in a linear manner that allows a cumulative BOLD signal derived from multiple individual events to be broken down into its constituent elements.
- (4) We illustrated event-related fMRI through a study of how individual visual experiences are encoded into memory.

### Neural Correlates of the BOLD Signal

- (1) Functional magnetic resonance imaging (fMRI) provides a measure of blood flow in terms of levels of blood oxygenation (the BOLD signal), giving an index of cognitive activity.
- (2) This raises the question of how this cognitive activity is related to neural activity.
- (3) One possibility is that cognitive activity detected by fMRI is correlated with the outputs of populations of neurons (as manifested in their firing activity). Another possibility is that the correlation is with the input to populations of neurons (as measured by the local field potential).
- (4) The experiments of Logothetis and his collaborators seem to show that the correlation is with the input to neural areas, rather than with their output.

## Further Reading

Ungerleider and Mishkin's paper "Two cortical visual systems" is reprinted in Cummins and Cummins 2000. Mishkin, Ungerleider, and Macko 1983/2001 is a little more accessible. David Milner and Melvyn Goodale have developed a different version of the two visual systems hypothesis, placing much more emphasis on studies of brain-damaged patients. See, for example, their book *The Visual Brain in Action* (2006). A more recent summary can be found in Milner and Goodale 2008 (including discussion of Ungerleider and Mishkin). A different development in terms of vision for action versus vision for higher mental processes has been proposed by the cognitive neuroscientist Marc Jeannerod, as presented in *Ways of Seeing*, coauthored with the philosopher Pierre Jacob (Jacob and Jeannerod 2003). A recent critique of the two-system account (with commentary from Milner, Goodale, and others) can be found in Schenk and McIntosh 2010. See Rossetti et al. 2017 for an up-to-date overview.



*The Handbook of Brain Theory and Neural Networks* (Arbib 2003) is the most comprehensive single-volume source for different types of computational neuroscience and neural computing, together with entries on neuroanatomy and many other “neural topics.” It contains useful introductory material and “road maps.” Dayan and Abbott 2005 and Trappenberg 2010 are other commonly used introductory textbooks. Scholarpedia.org is also a good source for introductory articles specifically on topics in computational neuroscience. McLeod, Plunkett, and Rolls 1998 is a good introduction to connectionism that comes with software allowing readers to get hands-on experience in connectionist modeling. Bechtel and Abrahamsen 2002 is also to be recommended. Useful article-length presentations are Rumelhart 1989 (in Posner 1989, reprinted in Haugeland 1997) and Churchland 1990b (in Cummins and Cummins 2000). A more recent discussion of connectionism can be found in McClelland et al. 2010, with commentary and target articles from others in the same issue. The mine/rock network described in the text was first presented in Gorman and Sejnowski 1988 and is discussed in Churchland 1990a. There are more references to literature on neural networks in the Further Reading sections for Chapters 5 and 10. Recent interest in neural networks has been associated with deep learning, discussed in Chapter 12.

A very readable book introducing PET and functional neuroimaging in general is Posner and Raichle 1994, written by two senior scientists participating in the lexical access experiments discussed in the text. These experiments are discussed in the article by Petersen et al. cited in the text and also (more accessibly) in Petersen and Fiez 2001. Rowe and Frackowiak 2003 is an article-length introduction to the basic principles of functional neuroimaging. Another good introduction to neuroimaging, including discussion of many of the experiments mentioned in this chapter (and with a lot of colorful illustrations), is Baars and Gage 2010.

For the specifics of event-related fMRI see the overviews in Buckner 1998 and Huettel 2012. The Huettel paper is published in a special issue of the journal *Neuroimage* edited by Peter Bendettini and titled “20 years of fMRI” (62:2, August 2012). The study referenced in the text is Brewer et al. 1998. See also Wagner et al. 1998 in the same issue of *Science*. For the first study using event-related fMRI, see Blamire et al. 1992. Dale and Buckner 1997 was another early study.

For specific references on the fMRI technology, see the suggestions for further reading in Chapter 9. For a survey of some of the general issues in thinking about the neural correlates of the BOLD signal, see Heeger and Ress 2002 and Raichle and Mintun 2006. Logothetis’s single-authored 2001 paper in the *Journal of Neuroscience* is a good introduction to the general issues as well as to his own experiments. See also Logothetis 2002. A more recent summary can be found in Goense, Whittingstall, and Logothetis 2012. For the Rees–Friston–Koch hypothesis, see Rees, Friston, and Koch 2000. For commentary on Logothetis, see Bandettini and Ungerleider 2001. For an alternative view, see Mukamel et al. 2005. Ekstrom 2010 discusses apparent dissociations between the BOLD signal and local field potentials.

## PART II

# MODELS AND TOOLS









## CHAPTER FOUR

# Physical Symbol Systems and the Language of Thought

### OVERVIEW 99

#### 4.1 The Physical Symbol System Hypothesis 100

Symbols and Symbol Systems 101  
Transforming Symbol Structures 102  
Intelligent Action and the Physical Symbol System 106

#### 4.2 From Physical Symbol Systems to the Language of Thought 106

Intentional Realism and Causation by Content 108

The Language of Thought and the Relation between Syntax and Semantics 110

#### 4.3 The Russian Room Argument and the Turing Test 114

Responding to the Russian Room Argument 117



## Overview

The analogy between minds and digital computers is one of the most powerful ideas in cognitive science. The physical symbol system hypothesis, proposed in 1975 by the computer scientists Herbert Simon and Allen Newell, articulates the analogy very clearly. It holds that all intelligent behavior essentially involves transforming physical symbols according to rules. Section 4.1 explains the basic idea, while Section 4.2 looks at the version of the physical symbol system hypothesis developed by the philosopher Jerry Fodor. Fodor develops a subtle and sophisticated argument for why symbolic information processing has to take place in a language of thought.

Both the general physical symbol system hypothesis and the language of thought hypothesis distinguish sharply between the syntax of information processing (the physical manipulation of symbol structures) and the semantics of information processing. The philosopher John Searle has developed a famous argument (the Chinese room argument) aiming to show that the project of modeling the mind as a computer is fatally flawed. We look at a version of his argument and at some of the ways of replying to it in Section 4.3.



## 4.1 The Physical Symbol System Hypothesis

In 1975 the Association of Computing Machinery gave their annual Turing Award to two very influential pioneers of artificial intelligence – Herbert Simon and Allen Newell. As a great example of the interdisciplinary nature of cognitive science, Simon was actually an economist and political scientist, rather than a computer scientist (as Newell was). Their joint contributions to computer science included the Logic Theory Machine (1957) and the General Problem Solver (1956), two early and very important programs that developed general strategies for solving formalized symbolic problems.

Newell and Simon gave a public lecture as one of the conditions of receiving the award. That lecture proposed a bold strategy both for the study of the human mind and for the emerging field of artificial intelligence (AI). Their manifesto hinged on what they called the *physical symbol system hypothesis*, which they proposed as a fundamental law for studying intelligence.

For Newell and Simon, the physical symbol system hypothesis is as basic to AI as the principle that the cell is the basic building block of all living organisms is to biology. Here is how they phrased it:

***The physical symbol system hypothesis:*** A physical symbol system has the necessary and sufficient means for general intelligent action.

There are two separate claims here. The first (the necessity claim) is that nothing can be capable of intelligent action unless it is a physical symbol system. So, since humans are capable of intelligent action, the human mind must be a physical symbol system. The second (the sufficiency claim) is that there is no obstacle in principle to constructing an artificial mind, provided that one tackles the problem by constructing a physical symbol system.

The significance of these two claims depends on what a physical symbol system is. Here are Newell and Simon again:

A physical symbol system consists of a set of entities, called symbols, which are physical patterns that can occur as components of another type of entity called an expression (or symbol structure). Thus a symbol structure is composed of a number of instances (or tokens) of symbols related in some physical way (such as one token being next to another). At any instant of time the system will contain a collection of these symbol structures. Besides these structures, the system also contains a collection of processes that operate on expressions to produce other expressions: processes of creation, modification, reproduction, and destruction. A physical symbol system is a machine that produces through time an evolving collection of symbol structures.

This passage illustrates four distinctive features of physical symbol systems. Here they are:

- 1 Symbols are physical patterns.
- 2 These symbols can be combined to form complex symbol structures.



- 3 The physical symbol system contains processes for manipulating symbols and complex symbol structures.
- 4 The processes for generating and transforming complex symbol structures can themselves be represented by symbols and symbol structures within the system.

You might have noticed that a physical symbol system looks very much like an abstract characterization of a digital computer. That is absolutely correct, as we'll now see.

## Symbols and Symbol Systems

To illustrate the first two ideas in the physical symbol system hypothesis we can go back to Turing machines, which we first encountered in Section 1.2 as abstract models of computation. Newell and Simon make clear in their paper how Turing's work on Turing machines in the 1930s was the first step toward the physical symbol system hypothesis.

(1) *Symbols are physical patterns.* For Newell and Simon symbols are physical objects, just as written letters on a page are physical objects, or spoken words (which are soundwaves). The symbols in Turing machines are also physical objects. They are inscriptions on the tape that the Turing machine is able to read. What the machine does at any given moment is fixed by the state it is in and the symbol that is on the cell being scanned.

Don't take this too literally, though. Even though a computer has an alphabet composed of the digits 0 and 1, we will not find any 0s and 1s in it if we open it up. If we dig down deep enough, all that there is to a computer is electricity flowing through circuits. If an electrical circuit functions as an on/off switch, then we can view that switch in symbolic terms as representing either a 0 (when it is off) or a 1 (when it is on). But there are no digits to be found in the circuit.

(2) *Symbols can be combined to form complex symbol structures.* Just as letters can be put together to form words, the symbols in any physical symbol system can be combined to form word-like symbol structures. Those word-like structures can then be put together to form sentence-like structures. Both types of combination are governed by strict rules. You can think of these strict rules as telling the symbol system which combinations of symbols count as grammatical.

These rules are likely to be *recursive* in form. That means that they will show how to get from an acceptable combination of symbols to a more complex combination that is still acceptable. The rules for how to define what counts as a sentence in the branch of logic known as sentential logic or propositional logic provide a good illustration of recursive rules and how they work. See Box 4.1.

Turing machines can scan only a single cell at a time, but they are still capable of working with complex symbol structures because those complex symbol structures can be built up from individual symbols in adjacent cells. The Turing machine needs to know two things: It needs to know what symbols can follow other symbols. And it needs some way of marking the end of complex symbols. The first can come from instructions in the machine table, while for the second there are symbols that serve as

### BOX 4.1 Defining Sentences in Propositional Logic

Propositional logic studies the logical relations holding between whole sentences, or propositions. The language of propositional logic is very simple. It contains basic symbols for sentences (such as "P," "Q," and "R"), together with a small set of logical connectives.

A typical formulation of propositional logic might have three connectives (the so-called Boolean connectives). These are " $\neg$ ," read as "not-"; " $\vee$ ," read as "or"; and " $\wedge$ ," read as "and."

These logical connectives allow sentence symbols to be combined to form more complex sentences. So, for example, " $P \wedge Q$ " is a sentence. It is true just when the two sentences P and Q are both true.

Propositional logic has clear and unambiguous rules for determining what counts as a legitimate sentence. These rules fix when the rules governing the connectives have been correctly applied. The legitimate combinations of symbols in the alphabet might typically be defined as follows.

- (a) Any sentence symbol is a sentence.
- (b) If " $\varphi$ " is a sentence then " $\neg \varphi$ " is a sentence.
- (c) If " $\varphi$ " and " $\psi$ " are sentences, then " $\varphi \wedge \psi$ " is a sentence.
- (d) If " $\varphi$ " and " $\psi$ " are sentences, then " $\varphi \vee \psi$ " is a sentence.

These are examples of what are called *recursive rules*. They show how, starting with a basic set of sentences (the sentence symbols), you can construct arbitrarily complex formulas that will count as genuine sentences.

Note that " $\varphi$ " and " $\psi$ " can stand here for any formula, not just for sentence symbols. So you can apply the recursive definition to show that  $\neg(P \wedge \neg P)$  is a genuine sentence of propositional logic.

Can you see how? (Hint: If P is a sentence symbol, then it is a sentence, by (a). If P is a sentence, then so is  $\neg P$ , by (b). Continue in this vein.)

punctuation marks, effectively telling the scanner when it has arrived at the end of a complex symbol.

## Transforming Symbol Structures

The third feature of physical symbol systems involves transformation.

(3) *The physical symbol system contains processes for manipulating symbols and complex symbol structures.* Here we have the distinctive claim of the physical symbol system hypothesis. Thinking is no more (and no less) than transforming symbol structures according to rules. Any system that can transform symbol structures in a sophisticated enough way will qualify as intelligent. According to Newell and Simon when we fully understand what is going on in intelligent agents (such as human beings), all we will ultimately find is symbol structures being transformed in rule-governed ways.



In the background here is Newell and Simon's fundamental idea that the essence of intelligent thinking is the ability to solve problems. Intelligence consists in the ability to work out, when confronted with a range of options, which of those options best matches certain requirements and constraints.

Intelligence cannot be applied without what might abstractly be called a search-space. The notion of a search-space is very general. Consider, for example, the position of one of the players halfway through a chess match. Each chess player has a large number of possible moves and a clearly defined aim – to checkmate her opponent. The possible moves define the search-space and the problem is deciding which of the possible moves will move her closest to her goal.

Another example (much studied by computer scientists and mathematicians) is a traveling salesperson who starts in a particular city (say, Boston) and has to visit twenty other cities as quickly and efficiently as possible before eventually returning to Boston. Here we can think about the search-space in terms of all the possible routes that start and end in Boston and go through the twenty cities (perhaps visiting some more than once). The diagram at the top in Figure 4.1 illustrates a simpler traveling salesperson problem with only five cities ( $a$ ,  $b$ ,  $c$ ,  $d$ , and  $e$ ).

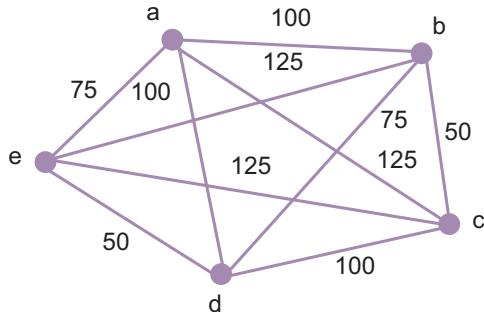
Search-spaces are typically represented in terms of states. There is an initial state (the start state) and a set of permissible transformations of that start state. The search-space consists of all the states that can be reached from the start state by applying the permissible transformations. The transformations can be carried out in any order. In the chess example, the start state is a particular configuration of the chess pieces and the permissible transformations are the legal moves in chess. In the traveling salesman problem, the start state might be Boston, for example, and the permissible transformations are given by all the ways of getting directly from one city to another. This means that each state of the traveling salesman problem is given by the current city, together with the cities already covered and the cities still left to visit.

Computer scientists standardly represent search-spaces in terms of trees. So, for example, the search-space for the traveling salesperson problem is given by a tree whose first node is the starting city. The diagram at the bottom of Figure 4.1 illustrates a part of the search-space for our five-city version of the traveling salesperson problem. A branch from the first node ( $a$ , the start city) goes to a node representing each city to which the start city is directly connected – i.e., cities  $b$ ,  $c$ ,  $d$ , and  $e$ . From each of those nodes, further branches connect each city to all the other cities to which it is directly connected. And so on.

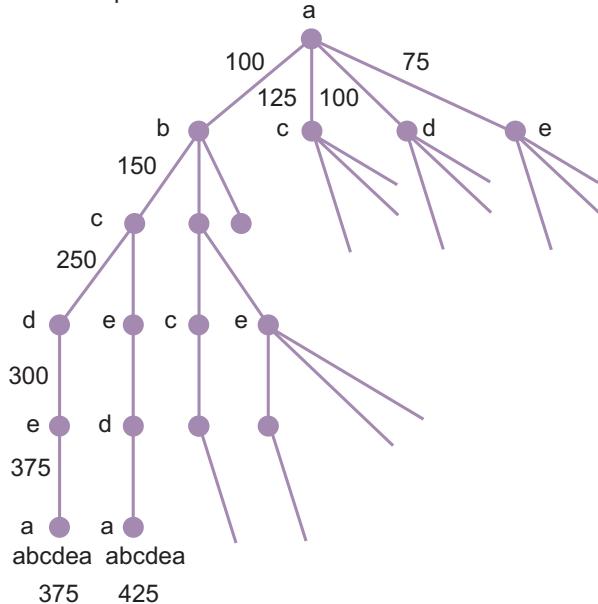
What counts as solving a problem? Basically, you've solved a problem when you've found the solution state in the search-space. In chess, the solution state is any configuration of the board in which the opponent's king is in checkmate. In the traveling salesperson problem, the solution is the shortest branch of the tree that ends with Boston and that has nodes on it corresponding to each city that the salesman needs to visit.

But how is that done? Obviously, you have to search through the search-space until you find a solution state. But this can be much harder than it sounds. Brute force searches that follow each branch of the tree typically only work for very simple problems. It does not

An instance of the traveling salesman problem



Search-space



**Figure 4.1** A typical traveling salesperson problem. The top diagram depicts the problem. A traveling salesperson has to find the shortest route between five cities. The diagram below depicts part of the search-space. A complete representation of the search-space would show twenty-four different routes.

take long for a problem space to get so big that it cannot be exhaustively searched in any feasible amount of time.

The traveling salesperson problem is a great example. If there are  $n$  cities, then it turns out that there are  $(n - 1)!$  possible routes to take into account, where  $(n - 1)! = (n - 1) \times (n - 2) \times (n - 3) \dots$ . This number of routes is not too many for the five-city version of the problem depicted in Figure 4.1 (it gives twenty-four different routes). But the problem gets out of control very quickly. A twenty-city version gives approximately  $6 \times 10^{16}$  different ways for a traveling salesperson to start in Boston and travel through the other nineteen cities visiting each exactly once. Checking one route per second, we would need more than the entire history of the universe to search the problem space exhaustively.



Newell and Simon developed their General Problem Solver (GPS) program as a way of solving problems of this type (although they focused on much simpler problems than the traveling salesperson problem, which still has no general solution).

The basic idea behind the GPS program is relatively straightforward. The program uses *means–end analysis*. Here are simple instructions for applying means–end analysis.

- 1 Evaluate the difference between the current state and the solution state.
- 2 Identify a transformation that reduces the difference between the current state and the solution state.
- 3 Check that the transformation in (2) can be applied to the current state.
  - 3a. If it can, then apply it and go back to step (1).
  - 3b. If it can't, then return to (2).

Means–end analysis is an example of what Newell and Simon call *heuristic search*. Heuristic search techniques are techniques for searching through a search-space that do not involve exhaustively tracing every branch in the tree until a solution is found. Heuristic search techniques reduce the size of the search-space in order to make the search process more manageable.



### Exercise 4.1 Explain how means–end analysis makes it more manageable to search the search-space.

Here is the problem of the foxes and the chickens – a type of problem that Newell and Simon showed could be solved by their GPS program. Imagine that there are three chickens and three foxes on one side of a river and they all need to get over to the other side. The only way to cross the river is in a boat that can take only two animals (or fewer) at a time. The boat can cross in either direction, but if at any point the foxes outnumber the chickens then the outnumbered chickens will be eaten. How can you get all the chickens and foxes onto the other side of the river without any of the chickens being eaten?

Here each state specifies which animals are on each bank and which in the boat (as well as the direction in which the boat is traveling). The start state obviously has all six on one bank (say the left bank) with nobody in the boat or on the other bank. The solution state is the state that has all six on the right bank, with nobody in the boat or on the other bank. The permissible transformations are defined by the rule that the boat cannot carry more than two animals at a time.

The foxes and chickens problem is a great example of how the GPS program works. If we feed into the GPS program representations of the start state and the solution, or goal state, the program employs various transformation strategies to minimize the difference between the start state and the goal state. The eventual solution is a series of representations, whose first member is a representation of the start state and whose final member is a representation of one of the goal states, and where each member is derived from its predecessor by a permissible transformation.



### Exercise 4.2 Find a solution to the foxes and chickens problem.

A final point. These rule-governed transformations are *algorithmic*. An algorithm is a finite set of unambiguous rules that can be applied systematically to transform an object or set of objects in definite and circumscribed ways. Algorithms are purely mechanical procedures. They can be followed blindly, without any exercise of judgment or intuition. Elementary school arithmetic provides plenty of examples of algorithms, such as the algorithms for multiplying pairs of numbers and for long division.

## Intelligent Action and the Physical Symbol System

The final feature of physical symbol systems is what really makes symbol systems capable of intelligent action.

(4) *The processes for generating and transforming complex symbol structures can themselves be represented by symbols and symbol structures within the system.* A fundamental feature of modern computers – so familiar that most of us never think about it – is the fact that a single computer (a single piece of hardware) can run many different programs, often simultaneously. This capability is what distinguishes a general-purpose computer from a specialized computing machine such as a pocket calculator. Computers can be programmed in this way because they can contain symbol structures that encode information about, and instructions for, other symbol structures.

Alan Turing proved that it is possible to construct a special kind of Turing machine (a *universal Turing machine*) that can mimic any specialized Turing machine implementing a particular algorithm. The universal Turing machine is a general-purpose computer. You can think of the specialized computers as software programs that run on the more general operating system of the universal Turing machine. The universal Turing machine is possible because Turing machine tables can be encoded as numbers, and hence can serve as inputs to Turing machines. The physical symbol system hypothesis builds something like this feature into the characterization of an intelligent system.

### 4.2

## From Physical Symbol Systems to the Language of Thought

The physical symbol system hypothesis tells us that intelligent agents solve problems by physically transforming symbolic structures. But we still need to know what these symbolic structures are, how they are transformed, and how those transformations give rise to intelligent action of the sort that human beings might carry out. This section looks at a proposal for answering these questions. This is the *language of thought hypothesis* developed by the philosopher and cognitive scientist Jerry Fodor (1935–2017).

According to Fodor's language of thought hypothesis, the basic symbol structures in the mind that carry information are sentences in an internal language of thought (sometimes called Mentalese). Information processing works by transforming those sentences in the language of thought.



Our starting point for exploring this idea is the basic fact that the mind *receives* information about its environment. Some of this information is carried by light waves arriving at the retina or sound waves hitting the eardrum. But in general, our behavior is not *determined* by the information that we receive. Different people, or the same person at different times, react differently to the same situation. There is no standard response to the pattern of sound waves associated (in English) with a cry of "Help!" for example. How we behave depends upon what our minds do with the information that they receive – how they *process* that information. If I run to your assistance when you cry "Help!" it is because my mind has somehow managed to decode your utterance as a word in English, worked out what you are trying to communicate, and then decided how to respond. This is all complex processing of the initial information that arrived at my eardrum.

But how does this information processing take place? How do vibrations on the ear drum lead to the muscle contractions involved when I save you from drowning? The information has to be carried by something. We know how the information is carried in the auditory system. We know that vibrations in the eardrum are transmitted by the ossicles to the inner ear, for example. What happens the further away the information travels from the eardrum is not so well understood, but another integral part of the general picture of the mind as physical symbol system is that there are physical structures that carry information and, by so doing, serve as *representations* of the immediate environment (or, of course, of things that are more abstract and/or more remote).

It is a basic assumption of cognitive science that information processing is, at bottom, a matter of transforming these representations in a way that finally yields the activity in the nervous system that "instructs" my limbs to jump into the water.

Information processing involves many different kinds of representation. This is illustrated by the example just given. The whole process begins with representations that carry information about vibrations in the eardrum. Somehow these representations get transformed into a much more complex representation that we might describe as my *belief* that you are in danger. This belief is in an important sense the "motor" of my behavior (my jumping into the water to rescue you). But it is not enough on its own. It needs to interact with other representations (such as my belief that I can reach you before you drown, and my *desire* to rescue you) in order to generate what I might think of as an *intention* to act in a certain way. This intention in turn gives rise to further representations, corresponding to the motor instructions that generate and control my bodily movements.

Among all these different types of representation, Fodor is particularly interested in the ones that correspond to beliefs, desires, and other similar psychological states. These psychological states are often called *propositional attitudes* by philosophers. They are called this because they can be analyzed as attitudes to propositions. Propositions are the sorts of thing that are expressed by ordinary sentences. So, there is a proposition expressed by the sentence "That person will drown" or by the sentence "It is snowing in St. Louis." Thinkers can have different attitudes to those propositions. I might fear the first, for example, and believe the second.

Fodor's starting point in thinking about propositional attitudes is that we are, by and large, pretty good at explaining and predicting other people's behavior in terms of what

they believe about the world and what they want to achieve. He thinks that this success is something that itself needs explanation. Why is the vocabulary of beliefs and desires (our *belief–desire psychology* or *propositional attitude psychology*) so deeply ingrained in us? Why does it seem so indispensable in our social interactions and social coordination? How and why do explanations that appeal to beliefs and desires actually work? And, in particular, why are these explanations so successful?

According to Fodor, there can be only one possible answer. Belief–desire psychology is successful because it is true. There really are such things as beliefs and desires. They are physical items that cause us to behave in certain ways. Belief–desire explanations are successful when they correctly identify the beliefs and other states that caused us to act in the way that we did.

If we say that someone jumped into the water because she believed that a child was drowning and wanted to save him, then what we are really claiming is that that person's bodily behavior was caused by internal items corresponding to the belief that someone is drowning and the desire to save her. This view is often called *intentional realism*.

Fodor's argument for the language of thought hypothesis is, in essence, that the hypothesis of intentional realism is the only way of explaining how belief–desire explanations can work. We will examine his argument in the next two subsections.



### **Exercise 4.3 Explain intentional realism in your own words.**

## **Intentional Realism and Causation by Content**

Intentional realism treats beliefs and desires as the sorts of things that can cause behavior. But this is a special type of causation. There is a fundamental difference between my leg moving because I am trying to achieve something (perhaps the journey of a thousand miles that starts with a single step) and my leg moving because a doctor has hit my knee with his hammer. In the first case, what causes my movement is what the desire is a desire for, namely, the beginning of the journey of a thousand miles. This is what philosophers call the *content* of the desire. There is nothing corresponding to a desire with content (or any other state with content) when a doctor hits my knee with a hammer. The movement that I make is simply a response to physical stimulus. It is not a response to something that I want to achieve.

This phenomenon, often called *causation by content*, is something that any version of intentional realism has to explain. That means taking into account the rational relations between belief and desires, on the one hand, and the behavior that they cause on the other. Beliefs and desires cause behavior that makes sense in light of them. Moving my leg is a rational thing to do if I desire to begin the journey of a thousand miles and believe that I am pointing in the right direction.

Yet, on the face of it, causation by content is deeply mysterious. It depends upon representations (stored information about the environment). In one sense representations are simply objects like any other – they might be patterns of sound waves, populations of neurons, or pieces of paper. Thought of in this way, it is no more difficult to understand



how representations can cause behavior than it is to understand how the doctor's hammer can make my leg move.

But the representations that we are interested in (such as beliefs and desires) are also things that bear a special *semantic* relation to the world – they have meanings. And the effects that representations have in the world is a function of what they mean. So, the puzzle, therefore, is how representations can have causal effects within the world as a function of their semantic properties, as a function of the relations in which they stand to other objects in the world (and indeed to objects that may not in fact even be in existence).

Fodor, along with almost all cognitive scientists and the vast majority of philosophers, holds that brains and the representations that they contain are physical entities. This means that brains can only be sensitive to certain types of property in mental representations. My utterance of the word "cat" is ultimately no more than a particular pattern of sound waves. These sound waves have certain physical properties (amplitude, frequency, wavelength, and so on) that can have certain effects on the brain. But the fact that those sound waves represent cats for English-speakers is a very different type of property (or at least, so the argument goes).

Let us call the physical properties that can be manipulated within brains *formal properties*. We call them this because they have to do with the physical *form* (i.e., the shape) of the representation. And let's use *semantics* for the properties that enable representations to represent – just as semantics is the branch of linguistics that deals with the meanings of words (how words represent).

This gives us another way of putting the problem. How can the brain be an information-processing machine if it is blind to the semantic properties of representations? How can the brain be an information-processing machine if all it can process are the formal properties of representations?



#### Exercise 4.4 Explain the contrast between formal and semantic properties in your own words.

At this point, we can see the particular slant that Fodor is putting on the physical symbol system hypothesis. Computers essentially manipulate strings of symbols. A computer programmed in binary, for example, manipulates strings of 1s and 0s. This string of 1s and 0s might represent a natural number, in the way that in binary 10 represents the number 2 and 11 represents the number 3. Or it might represent something completely different. It might represent whether or not the individual members of a long series of pixels are on or off, for example.

In fact, with a suitable coding, a string of 1s and 0s can represent just about anything. As far as the computer is concerned, however, what the string of 1s and 0s represents is completely irrelevant. The semantic properties of the string are irrelevant. The computer simply manipulates the formal properties of the string of 1s and 0s. In fact, it would be more accurate to say that the computer operates on *numerals* rather than *numbers*. Numerals are just symbols with particular shapes. Numbers are what those numerals represent.

But here's where the computer program comes in. The computer is programmed to manipulate strings of 1s and 0s in certain ways that yield the right result, even though the computer has no idea what that right result is. Take an adding machine, for example. Suppose it is given two strings of 0s and 1s and in response outputs a third string of 1s and 0s. If the first two strings represent the numbers 5 and 7, respectively, then (if the machine is well designed) the third string will be a binary representation of the number 12.

But even though all the computer is doing is mechanically manipulating 1s and 0s (numerals not numbers), operating on their formal properties, it nonetheless comes up with the right answer, all. So, although the computer itself is not concerned with what "12" means, the computer program must respect the rules of addition in order for the computational result – "12" – to have meaning.

In essence, what computer programmers do when they are programming an adding machine, is writing code so that purely mechanical manipulations of numerals will correctly track arithmetical relations between the numbers that the numerals represent. So, the adding machine must manipulate the numerals "7" and "5" in such a way that taking them as inputs to the machine yields the numeral "12," because it is an arithmetical fact that  $7 + 5 = 12$  (which is a statement about numbers, not numerals).

Fodor thinks that way of thinking about computer programs is a great model for the human brain. Brains are physical systems that can be sensitive only to the formal properties of mental representations. But nonetheless, as information-processing machines, they (like computers) have to respect the semantic properties of mental representations. The language of thought is what makes this possible.



#### **Exercise 4.5** Explain the analogy between brains and computers in your own words.

## The Language of Thought and the Relation between Syntax and Semantics

Here are the three main claims of Fodor's language of thought hypothesis.

- 1 Causation through content takes place through causal interactions between physical states.
- 2 These physical states have the structure of sentences, and their sentence-like structure determines how they are made up and how they interact with each other.
- 3 Causal transitions between sentences in the language of thought respect the rational relations between the contents of those sentences (what they mean).

According to Fodor, we think in sentences, but these are not sentences of a natural language such as English. The language of thought is much closer to a logical language, such as the propositional calculus (which we looked at briefly earlier in this chapter – see Box 4.1). It is supposed to be free of the nuances, ambiguities, and multiple layers of meaning that we find in English and other natural languages.



The analogy between the language of thought and logical languages is at the heart of Fodor's solution to the problem of causation by content. It is what lies behind claim (3). The basic fact about formal languages that Fodor exploits is the clear separation that they incorporate between *syntax* and *semantics*. Syntax has to do with symbols and the rules for combining them. You can think of it as the logical equivalent of grammar. Semantics, on the other hand, has to do with what the symbols actually mean and, relatedly, to what makes sentences true (or false).

To illustrate this general distinction, we can use the *predicate calculus*. This is a logical language more powerful and sophisticated than the propositional calculus we looked at in Box 4.1. Unlike the propositional calculus (which only allows us to talk about complete sentences or propositions) the predicate calculus allows us to talk directly about individuals and their properties. So, for example, the predicate calculus allows us to formalize inferences such as:

Hubert is laughing  
Therefore, someone is laughing

Or:

Everyone is laughing  
Therefore, Hubert is laughing

In order to represent these inferences, the predicate calculus has special symbols. These special symbols include individual constants that name particular objects, and predicate letters that serve to name properties. The symbols are typically identifiable by simple typographical features (such as uppercase for predicate letters and lowercase for individual constants) and they can be combined to make complex symbols according to certain rules. It also includes *quantifiers*, which are logical expressions corresponding to the English words "some" and "all."

From a syntactic point of view, a formal language such as the predicate calculus is simply a set of symbols of various types together with rules for manipulating those symbols according to their types. These rules identify the symbols only in terms of their typographical features. An example would be the rule that the space after an uppercase letter (e.g., the space in "F-") can be filled only with a lowercase letter (e.g., "a"). This rule is a way of capturing at the syntactic level the intuitive thought that properties apply primarily to things – because uppercase letters (such as "F-") can only be names of properties, while lowercase letters (such as "a") can only be names of objects. The rule achieves this, however, without explicitly stating anything about objects and properties. It just talks about symbols. It is a matter purely of the *syntax* of the language.

The connection between the formal system, on the one hand and what it is about, on the other, comes at the level of *semantics*. When we think about the semantics of a formal language, we assign objects to the individual constants and properties to the predicates. We identify the particular object that each individual constant names, for example. To provide a semantics for a language is to give an interpretation to the symbols it contains – to turn it from a collection of meaningless symbols into a representational system.



### Exercise 4.6 Explain the distinction between syntax and semantics in your own words.

Fodor's basic proposal is that we understand the relation between sentences in the language of thought and their content (or meaning) on the model of the relation between syntax and semantics in a formal system. Sentences in the language of thought can be viewed purely syntactically. From the syntactic point of view they are physical symbol structures composed of basic symbols arranged according to certain rules of composition. Or they can be viewed semantically in terms of how they represent the world (in which case they are being viewed as the vehicles of propositional attitudes).

So now, suppose we think that the causal transitions between sentences in the language of thought are essentially syntactic, that is, sensitive only to the formal properties of the relevant symbols, regardless of the symbols' meanings. Then we need to ask the following question:

Why do the syntactic relations between sentences in the language of thought map onto the semantic relations holding between the contents of those sentences?

If we take seriously the idea that the language of thought is a formal system, then this question has a perfectly straightforward answer. Syntactic transitions between sentences in the language of thought track semantic transitions between the contents of those sentences for precisely the same reason that syntax tracks semantics in any properly designed formal system.

Fodor can (and does) appeal to well-known results in meta-logic (the study of the expressive capacities and formal structure of logical systems). These results establish a significant degree of correspondence between syntax and semantics. So, for example, it is known that the first-order predicate calculus is sound and complete. That is to say, in every well-formed proof in the first-order predicate calculus the conclusion really is a logical consequence of the premises (*soundness*) and, conversely, for every argument in which the conclusion follows logically from the premises and both conclusion and premises are formulable in the first-order predicate calculus there is a well-formed proof (*completeness*).

The combination of soundness and completeness has the following important consequences. If a series of legitimate and formally definable syntactic transformations lead from formula A to a second formula B, then one can be sure that A cannot be true without B being true – and, conversely, if A entails B in a semantic sense then one can be sure that there will be a series of formally definable inferential transitions leading from A to B.

Here's an example. Suppose that we have two complex symbols, "Fa" and "Ga." Each of these symbols is a sentence in the language of thought with a particular syntactic shape. We know that "F–" and "G–" are symbols for *predicates*. Let us say that "F–" means "– is tall" and "G–" means "– has red hair." We also know that "a" is a *name* symbol. Let us say that "a" names Georgina. The meaning of "Fa" is that Georgina is tall, while the meaning of "Ga" is that Georgina has red hair.

Table 4.1 shows how a very simple piece of thinking might be analyzed by the language of thought hypothesis.

**TABLE 4.1** Syntax and semantics in the predicate calculus

SYMBOLS	TRANSFORMATION RULE	MEANING
1. $Fa$		1. Georgina is tall
2. $Ga$		2. Georgina has red hair
3. $(Fa \ \& \ Ga)$	If complex symbols "S" and "T" appear on earlier lines, then it is legitimate to write " $(S \ \& \ T)$ "	3. Georgina is tall and has red hair
4. $\exists x \ (Fx \ \& \ Gx)$	If on an earlier line there is a complex symbol containing a name symbol, then it is legitimate to replace the name symbol by " $x$ " and write " $\exists x -$ " in front of the complex symbol [NOTE: " $\exists x -$ " is the symbol for "there is at least one $x$ such that $-$ "]	4. At least one person is tall and has red hair

The table shows how two physical symbols: " $Fa$ " and " $Ga$ " are transformed in two inferential steps into the more complex physical symbol " $\exists x \ (Fx \ \& \ Gx)$ ." Here " $\exists x -$ " is the symbol for "there is at least one  $x$  such that  $-$ ," so that this sentence means "There is at least one thing that is both F and G."

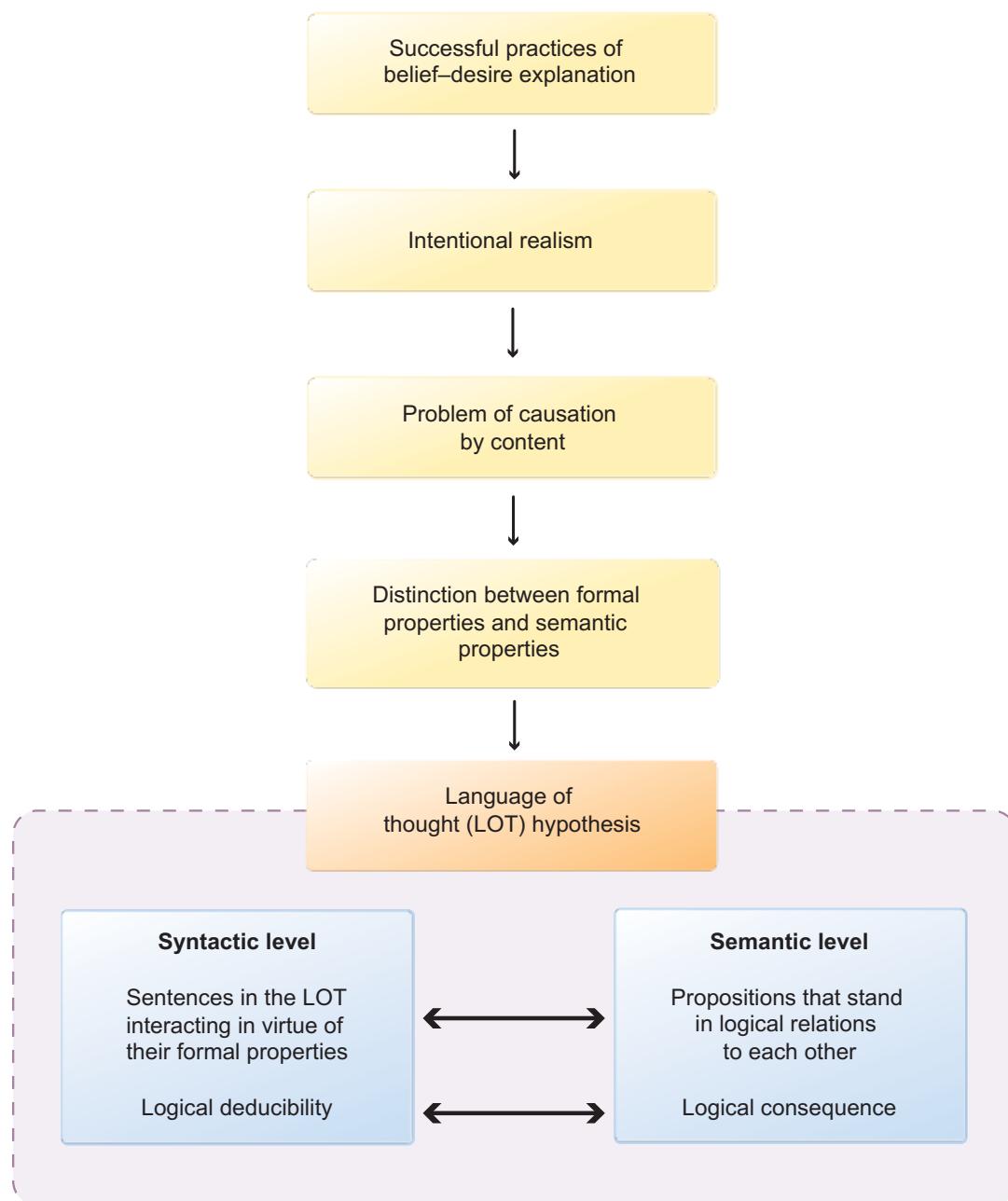
The rules that achieve this transformation are purely syntactic. They are simply rules for manipulating symbol structures. But when we look at the relation between the meanings of " $Fa$ " and " $Ga$ ," on the one hand, and the meaning of " $\exists x \ (Fx \ \& \ Gx)$ ," on the other, we see that those purely syntactic transformations preserve the logical relations between the propositions that the symbols stand for. If it is true that Georgina is tall and that Georgina has red hair, then it is certainly true that at least one person is tall and has red hair.

In sum, beliefs and desires are realized by language-like physical structures (sentences in the language of thought), and practical reasoning and other forms of thinking are ultimately just causal interactions between those structures. These causal interactions are sensitive only to the formal, syntactic properties of the physical structures. Yet, because the language of thought is a formal language with analogs of the formal properties of soundness and completeness, these purely syntactic transitions respect the semantic relations between the contents of the relevant beliefs and desires. This is how (Fodor claims) causation by content can take place in a purely physical system such as the human brain. And so, he argues, commonsense psychological explanation is vindicated by thinking of the mind as a computer processing sentences in the language of thought.

The line of reasoning that leads to the language of thought hypothesis is fairly complicated. To make it easier to keep track of the different steps, I have represented them diagrammatically in Figure 4.2.



**Exercise 4.7** Use the flowchart in Figure 4.2 to explain Fodor's argument in your own words.



**Figure 4.2** The structure of Fodor's argument for the language of thought hypothesis.

## 4.3

### The Russian Room Argument and the Turing Test

We need now to consider a fundamental objection to the very idea of the physical symbol system hypothesis. This objection comes from the philosopher John Searle, who is convinced that no machine built according to the physical symbol system hypothesis could possibly be capable of intelligent behavior.



Using a thought experiment, Searle tries to show that the physical symbol system hypothesis is completely mistaken. He rejects the idea that manipulating symbols is sufficient for intelligent behavior – even when the manipulation produces exactly the right outputs. What he tries to do is describe a situation in which symbols are correctly manipulated, but where there seems to be no genuine understanding and no genuine intelligence.

Searle asks us to imagine a machine that he calls a Russian room. (Searle originally used the example of Chinese for dramatic effect, since it is spoken and understood by very few people in the United States, and notoriously difficult to learn as a second language. I have decided to switch the example to Russian, for variety. Russian speakers should feel free to substitute their own example – an English room, for example, or an Arabic one.) Inside the room, a person receives pieces of paper through one window and passes out pieces of paper through another window. The pieces of paper have symbols in Russian written on them. The Russian room, in essence, is an input–output system, with symbols as inputs and outputs. The way the input–output system works is determined by a huge instruction manual that tells the person in the room which pieces of paper to pass out depending on which pieces of paper he receives.

The instruction manual is essentially just a way of pairing input symbols with output symbols. It is not written in Russian and can be understood and followed by someone who knows no Russian. All that the person needs to be able to do is to identify Russian symbols in some sort of syntactic way – according to their shape, for example. This is enough for them to be able to find the right output for each input – where the right output is taken to be the output dictated by the instruction manual.

The Russian room is admittedly a little far-fetched, but it does seem to be perfectly possible. Now, Searle continues, imagine two further things. Imagine, first, that the instruction manual has been written in such a way that the inputs are all questions in Russian and the outputs are all appropriate answers to those questions. To all intents and purposes, therefore, the Russian room is answering questions in Russian. Now imagine that the person in the room does not in fact know any Russian. All he is doing is following the instructions in the instruction manual (which is written in English) to match an outgoing piece of paper (the answer) to an incoming piece of paper (the question).

The situation is illustrated in Figure 4.3. What the Russian room shows, according to Searle, is that it is perfectly possible for there to be syntactic symbol manipulation without any form of intelligence or understanding. The Russian room as a whole does not understand Russian, because the man inside the room does not understand Russian.

The Russian room seems to be set up in accordance with the physical symbol system hypothesis. After all, the person in the Russian room is manipulating symbols according to their formal/syntactic properties. Moreover, the Russian room has been set up so that it produces the right output for every input. In the terms we used in the last section, the syntactic manipulation of the symbols preserves their semantic properties. The semantic properties of the input symbols are their meanings – i.e., certain questions. The semantic properties of the output symbols are answers to those questions. So, as long as the person in



**Figure 4.3** Inside and outside the Russian room.

the Russian room follows the instructions correctly, the semantic relations between input and output will be preserved. And yet, Searle argues, the Russian room does not understand Russian. How can it understand Russian, given that the person in the room does not understand Russian?

But if the Russian room does not understand Russian then, Searle argues, there is no sense in which it is behaving intelligently. To someone outside the room it might look as



if there is intelligent behavior going on. The machine does, after all, respond to the questions it is asked with answers that make sense. But this is just an illusion of intelligence. The Russian room cannot be behaving intelligently if it does not understand Russian. And so, it is a counterexample to the physical symbol system hypothesis – or so Searle argues.

In fact, Searle also thinks that the Russian room argument reveals a fundamental problem with the so-called Turing Test, proposed by Alan Turing in 1950 as a criterion for whether a machine was displaying real intelligence. Turing's basic idea is that, if an observer is communicating with a machine and cannot tell the difference between it and a human being, then that would show that the computer was genuinely intelligent. You might imagine, for example, that a judge is communicating simultaneously with the machine and with a human being via a computer screen and keyboards, but cannot tell which is which, however complicated and lengthy the interaction. For Turing, then, a machine that responds in exactly the way that a human being responds thereby counts as intelligent.



**Exercise 4.8** Explain in your own words why the Russian room argument is an objection to taking the Turing Test to reveal genuine intelligence.

## Responding to the Russian Room Argument

Many people have suggested that there is a crucial equivocation in the argument. The physical symbol system hypothesis is a hypothesis about how cognitive systems work. It says, in effect, that any cognitive system capable of intelligent behavior will be a physical symbol system – and hence that it will operate by manipulating physical symbol structures. The crucial step in the Russian room argument, however, is not a claim about the system as a whole. It is a claim about part of the system, namely, the person inside the room who is reading and applying the instruction manual. The force of the claim that the Russian room as a whole does not understand Russian rests almost entirely on the fact that this person does not understand Russian.

According to the *systems reply*, the Russian room argument is simply based on a mistake about where the intelligence is supposed to be located. Supporters of the systems reply hold that the Russian room as a whole understands Russian and is displaying intelligent behavior, even though the person inside the room does not understand Russian.

Searle himself is not very impressed by the systems reply. He has a clever objection. Instead of imagining yourself in the Russian room, imagine the Russian room inside you! If you memorize the instruction manual then, Searle says, you have effectively internalized the Russian room. Of course, it's hard to imagine that anyone could have a good enough memory to do this, but there are no reasons to think that it is in principle impossible. But, Searle argues, internalizing the Russian room in this way is not enough to turn you from someone who does not understand Russian into someone who does. After all, what you've memorized is not Russian, but just a complex set

of rules for mapping some symbols you don't understand onto other symbols you don't understand.



### Exercise 4.9 How might a defender of the systems reply to Searle's response?

Another common way of responding to the Russian room argument is the *robot reply*. According to the robot reply, Searle is right about the Russian room not understanding Russian, but wrong about the reasons why. The problem with the Russian room has nothing to do with some sort of impassable gap between syntax and semantics. The problem, rather, is that it is embodied agents who understand Russian, not disembodied cognitive systems into which pieces of paper enter and from which other pieces of paper come out. Understanding Russian is a complex ability that manifests itself in how an agent interacts with other people and with items in the world.

The ability to understand Russian involves, at a minimum, being able to carry out instructions given in Russian, to coordinate with other Russian speakers, to read Russian characters, and to carry on a conversation. In order to build a machine that could do all this we would need to embed the Russian room in a robot, providing it with some analog of sensory organs, vocal apparatus, and limbs. If the Russian room had all this and could behave in the way that a Russian-speaker behaves then, a supporter of the robot reply would say, there is no reason to deny that the system understands Russian and is behaving intelligently.

Again, Searle is unconvinced. For him the gulf between syntax and semantics can't simply be overcome by turning the Russian room into a Russian robot. An embodied Russian room might indeed stop when it "sees" the Russian word for "stop." But this would simply be something it has learned to do. It no more understands what the character means than a laboratory pigeon that has been trained not to peck at a piece of card with the same character on it. Interacting with the environment is not the same as understanding it. Even if the Russian robot does and says all the right things, this does not show that it understands Russian. The basic problem still remains, as far as Searle is concerned: simply manipulating symbols cannot create meaning and unless the symbols are meaningful to the Russian room there is no relation between what it does and what a "real" Russian-speaker might do.



### Exercise 4.10 Explain the robot reply in your own words and assess Searle's response to it.

The Russian room argument raises fundamental questions about the nature of intelligence and, relatedly, of what counts as genuine thinking. Searle often presents it as an objection to what is often called the project of *strong AI* – this is the project of building machines that are genuinely intelligent (as opposed to the project of *weak AI*, which simply aims to build machines that do the things that human beings can do). And this is what many supporters (and critics) of the Russian room argument have focused on.

It is important to realize, though, that, if the argument is sound, it strikes at the very possibility of the physical symbol system hypothesis. What Searle is trying to show



is that there must be more to genuine thinking than simply manipulating symbols according to rules, whereas the physical symbol system hypothesis says that that is all that thinking is. For that reason, the Russian room argument is a very useful tool for thinking about some of the broader, theoretical issues that the physical symbol system hypothesis raises.



## Summary

This chapter has looked at the physical symbol system hypothesis, originally proposed by Newell and Simon. This hypothesis says that thinking consists in manipulating symbol structures according to rules. After introducing the physical symbol system hypothesis, we considered Jerry Fodor's suggestion that these symbol structures are sentences in an internal language of thought. The chapter ended with the Russian room argument, an objection to the basic idea of the physical symbol system hypothesis.

## Checklist

**The physical symbol system hypothesis states that a physical symbol system has the necessary and sufficient means for general intelligent action. In more detail:**

- (1) The symbols are physical patterns.
- (2) Physical symbols can be combined to form complex symbol structures.
- (3) Physical symbol systems contain processes for manipulating complex symbol structures.
- (4) The processes for manipulating complex symbol structures can be represented by symbols and structures within the system.
- (5) Problems are solved by generating and modifying symbol structures until a solution structure is reached.

**The physical symbol system hypothesis is very programmatic. Fodor's language of thought hypothesis is one way of turning the physical symbol system hypothesis into a concrete proposal about how the mind works.**

- (1) The language of thought hypothesis is grounded in intentional realism. Psychological states such as belief and desire are real physical entities. These entities are sentences in the language of thought.
- (2) The hypothesis offers a way of explaining causation by content (i.e., how physical representations can have causal effects in the world as a function of how they represent the world).
- (3) Fodor suggests that we understand the relation between sentences in the language of thought and their contents on the model of the relation between syntax and semantics in a formal system.
- (4) The syntax of the language of thought tracks its semantics because the language of thought is a formal language with analogs of the formal properties of soundness and completeness.

The Russian room argument is a thought experiment directed against the idea that the rule-governed manipulation of symbols is sufficient to produce intelligent behavior.

- (1) The person in the Russian room is manipulating symbols according to their formal/syntactic properties without any understanding of Russian.
- (2) According to the systems reply, the Russian room argument misses the point, because the real question is whether the system as a whole understands Russian, not whether the person in the room understands Russian.
- (3) According to the robot reply, the Russian room does not understand Russian. But this is not because of any uncrossable gap between syntax and semantics. Rather, it is because the Russian room has no opportunity to interact with the environment and other people.

## Further Reading

The paper by Newell and Simon discussed in Section 4.1 is reprinted in a number of places, including Boden 1990b and Bermúdez 2006. A good introduction to the general ideas behind the physical symbol system hypothesis in the context of artificial intelligence is Haugeland 1985, particularly chapter 2, and Haugeland 1997, chapter 4. See also chapters 1–3 of Johnson-Laird 1988, chapters 4 and 5 of Copeland 1993, chapter 2 of Dawson 1998, and the *Encyclopedia of Cognitive Science* entry on symbol systems (Nadel 2005). Russell and Norvig 2009 is the third edition of a popular AI textbook. Also see Poole and Mackworth 2010, Warwick 2012, and Proudfoot and Copeland's chapter on artificial intelligence in *The Oxford Handbook of Philosophy of Cognitive Science* (Margolis, Samuels, and Stich 2012).

Fodor 1975 and 1987 are classic expositions of the language of thought approach from a philosophical perspective. For Fodor's most recent views, see Fodor 2008. For a psychologist's perspective, see Pylyshyn's book *Computation and Cognition* (1984) and his earlier target article in *Behavioral and Brain Sciences* (1980). More recent philosophical discussions of the language of thought can be found in Schneider 2011 and Schneider and Katz 2012. The *Encyclopedia of Cognitive Science* has an entry on the topic, as does the *Stanford Encyclopedia of Philosophy*. For a general, philosophical discussion of the computational picture of the mind, Crane 2003 and Sterelny 1990 are recommended. Block 1995a explores the metaphor of the mind as the software of the brain. Fodor's argument for the language of thought hypothesis is closely tied to important research in mathematical logic and the theory of computation. Rogers 1971 is an accessible overview. For general introductions to philosophical debates about mental causation and the more general mind–body problem, see Heil 2004 and Searle 2004.

Searle presents what I am calling the Russian room argument and he termed the Chinese room argument in his "Minds, brains, and programs" (1980). This was originally published in the journal *Behavioral and Brain Sciences* with extensive commentary from many cognitive scientists. A related problem, the symbol-grounding problem, is introduced and discussed in Harnad 1990 (available in the online resources). Margaret Boden's article "Escaping from the Chinese room" (1990a), reprinted in Heil 2004, is a good place to start in thinking about the argument. In January 1990, the periodical *Scientific American* devoted a special issue to the tenth anniversary



of Searle's argument. Churchland and Churchland 1990 was one of the contributions, arguing that while the argument is effective against classical AI, it leaves artificial neural networks untouched. Preston and Bishop 2002 is a collection of articles dedicated to the Russian room argument, covering the principal lines of response. The entry on the Chinese room argument in the online *Stanford Encyclopedia of Philosophy* is comprehensive and has a very full bibliography.





## CHAPTER FIVE

# Neural Networks and Distributed Information Processing

### OVERVIEW 123

- 5.1 Neurally Inspired Models of Information Processing** 124
  - Neurons and Network Units 125
- 5.2 Single-Layer Networks and Boolean Functions** 128
  - Learning in Single-Layer Networks: The Perceptron Convergence Rule 131
  - Linear Separability and the Limits of Perceptron Convergence 134

### 5.3 Multilayer Networks 137

- The Backpropagation Algorithm 138
- How Biologically Plausible Are Neural Networks? 139

### 5.4 Information Processing in Neural Networks: Key Features 141

- Distributed Representations 141
- No Clear Distinction between Information Storage and Information Processing 142
- The Ability to Learn from "Experience" 143



## Overview

This chapter looks at a very different approach to information processing. Neural networks are based on an idealized model of how neurons work. The chapter begins in Section 5.1 by reviewing some of the motivations for neurally inspired models of information processing and looking at how the individual units in neural networks compare to biological neurons.

The simplest artificial neural networks are single-layer networks. These are explored in Section 5.2. We will see that any digital computer can be simulated by a suitably chained together set of single-layer networks. However, they are limited in what they can learn.

Overcoming those limits requires moving from single-layer networks to multilayer networks, which are capable of learning through the backpropagation of error. In Section 5.3 we look at the backpropagation algorithm used to train multilayer networks. Finally, Section 5.4 summarizes the key features of information processing in multilayer artificial neural networks, explaining key differences between neural networks and physical symbol systems.

 5.1

## Neurally Inspired Models of Information Processing

We saw in Part I (particularly in Chapter 3) that detailed knowledge of how the brain works has increased dramatically in recent years. Neuroimaging techniques, such as fMRI and PET, have allowed neuroscientists to begin establishing large-scale correlations between types of cognitive functioning and specific brain areas. Combining this with the information available from studies of brain-damaged patients allows cognitive scientists to build up a functional map of the brain.

Other techniques have made it possible to study brain activity (in nonhuman animals, from monkeys to sea-slugs) at the level of the single neuron. Microelectrodes can be used to record electrical activity both inside a single neuron and in the vicinity of that neuron. Recording from inside neurons allows a picture to be built up of the different types of input to the neuron, both excitatory and inhibitory, and of the mechanisms that modulate output signals.

But none of these techniques offers direct insight into how information is processed in the brain. PET and fMRI are good sources of information about which brain areas are involved in particular cognitive tasks, but they do not tell us anything about how those cognitive tasks are actually carried out. We need to know not just *what* particular regions of the brain do, but *how* they do it. Nor will this information come from single-neuron recordings. We may well find out from single-neuron recordings in monkeys that particular types of neuron in particular areas of the brain respond very selectively to a narrow range of visual stimuli, but we have as yet no idea how to scale this up into an account of how vision works.

The brain is an extraordinarily complicated set of interlocking and interconnected circuits. The most fundamental feature of the brain is its *connectivity* and the crucial question in understanding the brain is how distributed patterns of activation across populations of neurons can give rise to perception, memory, sensorimotor control, and high-level cognition. But we have (as yet) limited tools for directly studying how populations of neurons work.

Since we do not have the equipment and resources to study populations of neurons directly, many researchers have developed techniques for studying populations of neurons indirectly. The new strategy is to construct models that approximate populations of neurons in certain important respects. These are called neural network models, or artificial neural networks.

There are many different types of neural network models and many different ways of using them. The focus in *computational neuroscience* is on modeling biological neurons and populations of neurons. Computational neuroscientists start from what is known about the biology of the brain and then construct models by abstracting away from some biological details while preserving others. *Connectionist modelers* often pay less attention to the constraints of biology. They tend to start with generic models. Their aim is to show how those models can be modified and adapted to simulate and reproduce well-documented psychological phenomena, such as the patterns of development that children



go through when they acquire language, or the way in which cognitive processes break down in brain-damaged patients.

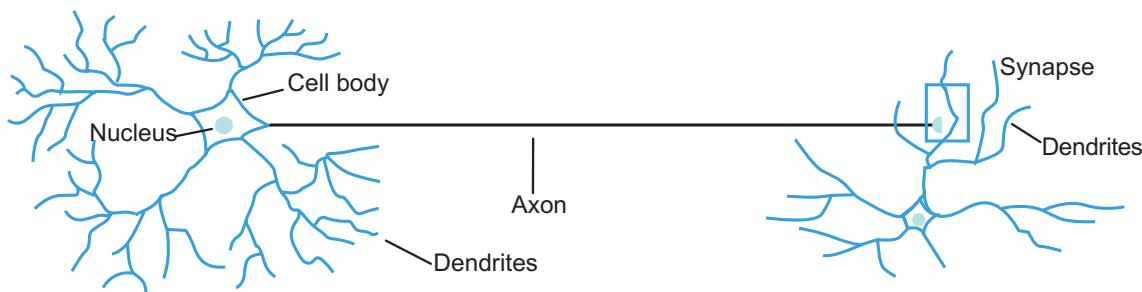
For our purposes here, the differences between computational neuroscientists and connectionist modelers are less important than what they have in common. Neural network models are distinctive in how they store information, how they retrieve it, and how they process it. And even those models that are not biologically driven remain neurally inspired. This neurally inspired way of thinking about information processing is the focus of this chapter.

## Neurons and Network Units

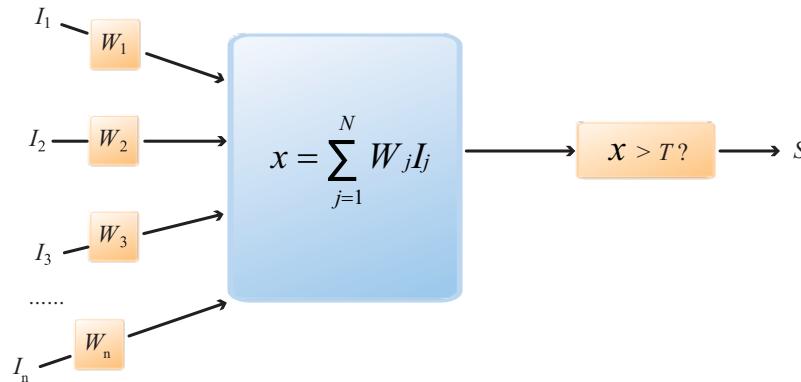
Neural networks are made up of individual units loosely based on biological neurons. There are many different types of neuron in the nervous system, but they all share a common basic structure. Each neuron is a cell and so has a cell body (a *soma*) containing a nucleus. There are many root-like extensions from the cell body. These are called *neurites*. There are two different types of neurite. Each neuron has many dendrites and a single axon. The dendrites are thinner than the axon and form what looks like a little bush (as illustrated in Figure 5.1). The axon itself eventually splits into a number of branches, each terminating in a little *endbulb* that comes close to the dendrites of another neuron.

Neurons receive signals from other neurons. A typical neuron might receive inputs from 10,000 neurons, but the number is as great as 50,000 for some neurons in the brain area called the hippocampus. These signals are received through the dendrites, which can be thought of as the receiving end of the neuron. A sending neuron transmits a signal along its axon to a *synapse*, which is the site where the end of an axon branch comes close to a dendrite or the cell body of another neuron. When the signal from the sending (or *presynaptic*) neuron reaches the synapse, it generates an electrical signal in the dendrites of the receiving (or *postsynaptic*) neuron.

The basic activity of a neuron is to fire an electrical impulse along its axon. The single most important fact about the firing of neurons is that it depends upon activity at the synapses. Some of the signals reaching the neuron's dendrites promote firing and others



**Figure 5.1** Schematic illustration of a typical neuron.



**Figure 5.2** An artificial neuron.

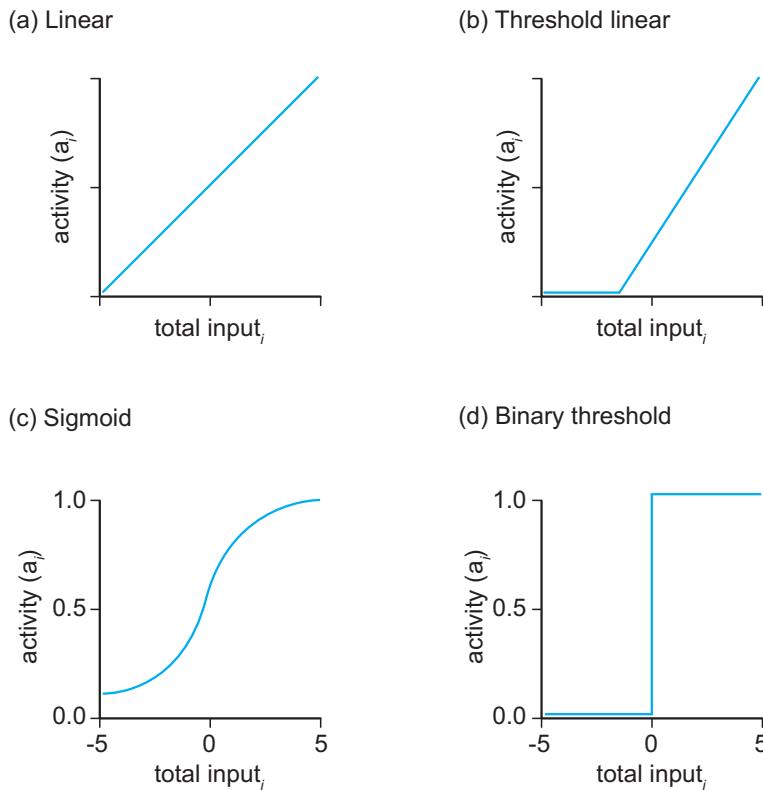
inhibit it. These are called *excitatory and inhibitory* synapses, respectively. If the sum of excitatory and inhibitory synapses exceeds the *threshold of the neuron* then the neuron will fire.

Neural networks are built up of interconnected populations of units that share some of the generic characteristics of biological neurons. Figure 5.2 illustrates a typical network unit. There are  $n$  inputs, corresponding to synaptic connections to presynaptic neurons. Signals from the presynaptic neurons might be excitatory or inhibitory. This is captured in the model by assigning a numerical weight  $W_i$  to each input  $I_i$ . Typically, the weight will be a real number between 1 and -1. A positive weight corresponds to an excitatory synapse and a negative weight to an inhibitory synapse.

Multiplying each input by its weight indicates the strength of the signal at each synapse. Adding all these individual *activation levels* together gives the total input to the unit, corresponding to the total signal reaching the nucleus of the neuron. This is represented using standard mathematical format in Figure 5.2. (A reminder –  $\Sigma$  is the symbol for summation [repeated addition]). The  $N$  above the summation sign indicates that there are  $N$  many things to add together. Each of the things added together is the product of  $I_j$  and  $W_j$  for some value of  $j$  between 1 and  $N$ .) If the total input exceeds the threshold ( $T$ ) then the neuron “fires” and transmits an output signal.

The one thing that remains to be specified is the strength of the output signal. We know that the unit will transmit a signal if the total input exceeds its designated threshold, but we do not yet know what that signal is. For this we need to specify an *activation function* – a function that assigns an output signal on the basis of the total input. Neural network designers standardly choose from several different types of activation function. Some of these are illustrated in Figure 5.3.

The simplest activation function is a linear function on which the output signal increases in direct proportion to the total input. (Linear functions are so called because they take a straight line when drawn on a graph.) The threshold linear function is a slight



**Figure 5.3** Four different activation functions. Each one fixes a neuron's activation level as a function of the total input to the neuron. (Adapted from McLeod, Plunkett, and Rolls 1998)

modification of this. This function yields no output signal until the total input reaches the threshold – and then the strength of the output signal increases proportionately to the total input. There is also a binary threshold function, which effectively operates like an on/off switch. It either yields zero output (when the input signal is below threshold) or maximum output (when the input signal is at or above threshold).

The threshold functions are intended to reflect a very basic property of biological neurons, which is that they only fire when their total input is suitably strong. The binary threshold activation function models neurons that either fire or don't fire, while the threshold linear function models neurons whose firing rate increases in proportion to the total input once the threshold has been reached.

The sigmoid function is a very commonly used nonlinear activation function. This reflects some of the properties of real neurons in that it effectively has a threshold below which total input has little effect and a ceiling above which the output remains more or less constant despite increases in total input. The ceiling corresponds to the maximum firing rate of the neuron. Between the threshold and the ceiling the strength of the output signal is roughly proportionate to the total input and so looks linear. But the function as a whole is nonlinear and drawn with a curve.

We see, then, how each individual unit in a network functions. The next step is to see how they can be used to process information. We will start out by looking at the simplest neural networks. These are *single-layer networks*.

## 5.2

## Single-Layer Networks and Boolean Functions

The first neural networks were studied in the 1940s and 1950s, pioneered by the neuroscientist Warren McCulloch and the logician Walter Pitts. They were known as single-layer networks. To see what single-layer networks can do (and what they can't do) we need so start with a quick refresher on *mapping functions*.

The basic idea of a function should be familiar, even if the terminology may not be. Addition is a function. Given two numbers as *inputs*, the addition function yields a third number as *output*. The output is the sum of the two inputs. Multiplication is also a function. Here the third number is the product of the two inputs.

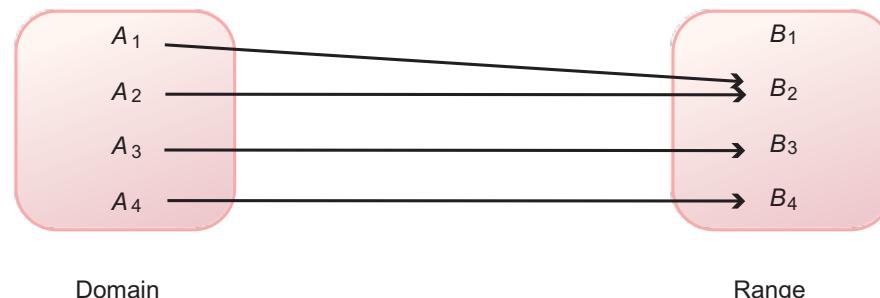
Some helpful terminology. Suppose that we have a set of items. We can call that a *domain*. Let there be another set of items, which we can call the *range*. A mapping function maps each item from the domain onto exactly one item from the range. The defining feature of a function is that no item in the domain gets mapped to more than one item in the range. Functions are *single-valued*. The operation of taking square roots, for example, is not a function (at least when negative numbers are included), since every positive number has two square roots.



**Exercise 5.1** Give another example of an arithmetical operation that counts as a function. And another example of an operation that is not a function.

Figure 5.4 gives an example of a mapping function. The arrows indicate which item in the domain is mapped to each item in the range. It is perfectly acceptable for two or more items in the domain to be mapped to a single item in the range (as is the case with  $A_1$  and  $A_2$ ). But, because functions are single-valued, no item in the domain can be mapped onto more than one item in the range.

Let's turn now to mapping functions of a very special kind. These are functions with a range consisting of two items, one corresponding to TRUE and the other corresponding to FALSE. We



**Figure 5.4** Illustration of a mapping function. A mapping function maps each item in its domain to exactly one item in its range.



can think about such functions as ways of classifying objects in the domain of the function. Imagine that the domain of the function contains all the natural numbers and the range of the function contains two items corresponding to TRUE and FALSE. Then we can identify any subset we please of the natural numbers by mapping the members of that subset onto TRUE and all the others onto FALSE. If the subset that the function maps onto TRUE contains all and only the even numbers, for example, then we have a way of picking out the set of the even numbers.

Now, we have all the machinery we need to introduce the so-called *binary Boolean functions*. These functions all have the same range as our even number function, namely, the set consisting of the two truth values TRUE and FALSE. Instead of having numbers in the domain, however, the domain of these functions is made up of pairs of truth values.

There are four different possible pairs of truth values. These pairs form the domain of the binary Boolean functions. The range, as with all Boolean functions, is given by the set {TRUE, FALSE}, as illustrated below:

DOMAIN	RANGE
FALSE, FALSE	
FALSE, TRUE	FALSE
TRUE, FALSE	TRUE
TRUE, TRUE	

Each binary Boolean function assigns either TRUE or FALSE to each pair of truth values.

You can think of a binary Boolean function as a way of showing how to fix the truth value of a complex sentence built up from two simpler sentences on the basis of the truth values of those simpler sentences. Some Boolean functions should be very familiar. There is a binary Boolean function standardly known as AND, for example. AND maps the pair {TRUE, TRUE} to TRUE and maps all other pairs of truth values to FALSE. To put it another way, if you are given a sentence A and a sentence B, then the only circumstance in which it is true to claim A AND B is the circumstance in which both A and B have the value TRUE.

Similarly, OR is the name of the Boolean function that maps the pair {FALSE, FALSE} to FALSE, and the other three pairs to TRUE. Alternatively, if you are given sentences A and B then the only circumstance in which it is false to claim A OR B is the circumstance in which both A and B have the value FALSE.

It is important that the OR function assigns TRUE to the pair {TRUE, TRUE}, so that A OR B is true in the case where both A and B are true. As we shall see, there is a Boolean function that behaves just like OR, except that it assigns FALSE to {TRUE, TRUE}. This is the so-called XOR function (an abbreviation of exclusive-OR). XOR *cannot* be represented by a single-layer network. We will look at this in more detail in Section 5.2.

We can represent these functions using what logicians call a truth table. The truth table for AND tells us how the truth value of A AND B varies according to the truth value of A and B, respectively (or, as a logician would say, as a *function* of the truth values of A and B).

A	B	A AND B
FALSE	FALSE	FALSE
FALSE	TRUE	FALSE
TRUE	FALSE	FALSE
TRUE	TRUE	TRUE



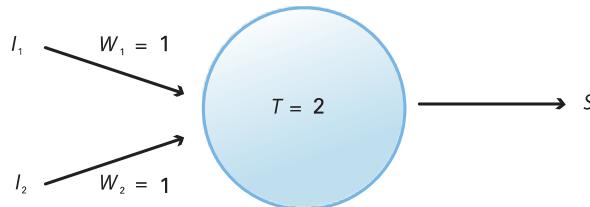
### Exercise 5.2 Give a truth table for the Boolean function OR.

Now, back to neural networks. The key point (first pointed out by McCulloch and Pitts) is that we can use very simple neural networks to represent some of the binary Boolean functions. The first step is to represent Boolean functions using numbers (since we need numbers as inputs and outputs for the arithmetic of the activation function to work). This is easy. We can represent TRUE by the number 1 and FALSE by 0, as is standard in logic and computer science. If we design our network unit so that it only takes 1 and 0 as inputs and only produces 1 and 0 as outputs, then it will be computing a Boolean function.

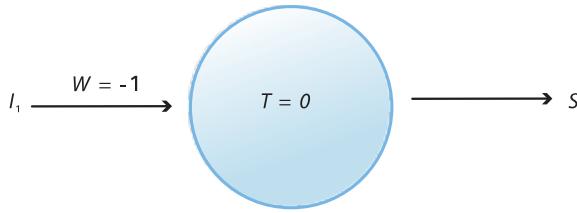
It is easy to see how we design our network unit to take only 0 and 1 as input. But how do we design it to produce only 0 and 1 as output?

Think back to the activation functions depicted in Figure 5.3, particularly binary threshold activation functions. These functions output 0 until the threshold is reached. Once the threshold is reached they output 1, irrespective of how the input increases. So, to represent a particular Boolean function, we need to set the weights and the threshold so that the network mimics the truth table for that Boolean function. A network that represents AND, for example, will have to output a 0 whenever the input is either (0, 0), (0, 1), or (1, 0). And it will have to output a 1 whenever the input is (1, 1).

The trick in getting a network to do this is to set the weights and the threshold appropriately. Look at Figure 5.5. If we set the weights at 1 for both inputs and the threshold at 2, then the unit will only fire when both inputs are 1. If both inputs are 1 then the total input is  $(I_1 \times W_1) + (I_2 \times W_2) = (1 \times 1) + (1 \times 1) = 2$ , which is the threshold. Since the network is using a binary threshold activation function (as described in the previous paragraph), in this case the output will be 1. If either input is a 0 (or both are) then the threshold will not be met, and so the output is 0. If we take 1 to represent TRUE and 0 to represent FALSE, then this network represents the AND function. It functions as what computer scientists call an AND-gate.



**Figure 5.5** A single-layer network representing the Boolean function AND.



**Figure 5.6** A single-layer network representing the Boolean function NOT.



### Exercise 5.3 Show how a network unit can represent OR and hence function as an OR-gate.

There are Boolean functions besides the binary ones. In fact, there are  $n$ -ary Boolean functions for every natural number  $n$  (including 0). But cognitive scientists are generally only interested in one nonbinary Boolean function. This is the unary function NOT. As its name suggests, NOT A is true if A is false and NOT A is false if A is true. Again, this is easily represented by a single network unit, as illustrated in Figure 5.6. The trick is to set the weights and threshold to get the desired result.



### Exercise 5.4 Explain why the network unit in Figure 5.6 represents the unary Boolean function NOT.

We see, then, single-layer networks can achieve a lot. As any computer scientist knows, modern digital computers are in the last analysis just incredibly complicated systems of AND-gates, OR-gates, and NOT-gates. So, by chaining together individual network units into a network we can do anything that can be done by a digital computer. (This is why I earlier said that cognitive scientists are generally only interested in one nonbinary Boolean function. AND, NOT, OR, and a little ingenuity are enough to simulate any  $n$ -ary Boolean function, no matter how complicated.)

There is something missing, however. As we have seen, the key to getting single units to represent Boolean functions such as NOT and OR lies in setting the weights and the threshold. But this raises some fundamental questions: How do the weights get set? How does the threshold get set? Is there any room for learning?

## Learning in Single-Layer Networks: The Perceptron Convergence Rule

In 1949 Donald Hebb published *The Organization of Behavior* in which he speculated about how learning might take place in the brain. His basic idea (the idea behind what we now call *Hebbian learning*) is that learning is at bottom an associative process. He famously wrote:

When an axon of a cell A is near enough to excite cell B or repeatedly or persistently takes part in firing it, some growth or metabolic change takes place in both cells such that A's efficiency, as one of the cells firing B, is increased.

Hebbian learning proceeds by synaptic modification. If A is a presynaptic neuron and B a postsynaptic neuron, then every time that B fires after A fires increases the probability that B will fire after A fires (this is what Hebb means by an increase in A's efficiency).

In slogan form, Hebbian learning is the principle that *neurons that fire together, wire together*. It has proved to be a very useful tool in modeling basic pattern recognition and pattern completion, as well as featuring in more complicated learning algorithms, such as the competitive learning algorithm discussed in Section 5.3.

Hebb was speculating about real neurons, not artificial ones. And, although there is strong evidence that Hebbian learning does take place in the nervous system, the first significant research on learning in artificial neural networks modified the Hebbian model very significantly. In the 1950s Frank Rosenblatt studied learning in single-layer networks. In an influential article in 1958 he called these networks *perceptrons*.

Rosenblatt was looking for a learning rule that would allow a network with random weights and a random threshold to settle on a configuration of weights and thresholds that would allow it to solve a given problem. Solving a given problem means producing the right output for every input.

The learning in this case is *supervised* learning. This means that, whenever the network produces the wrong output for a given input, it is told that it has made an error. The process of learning (for a neural network) is the process of changing the weights and/or the threshold in response to error. Learning is successful when these changes in the weights and/or the threshold converge upon a configuration that always produces the desired output for a given input.

Rosenblatt called his learning rule the *perceptron convergence rule*. The perceptron convergence rule has some similarities with Hebbian learning. Like Hebbian learning it relies on the basic principle that changes in weight are determined solely by what happens locally – that is, by what happens at the input and what happens at the output. But, unlike Hebbian learning, it is a supervised algorithm – it requires feedback about incorrect solutions to the problem the network is trying to solve.

The perceptron convergence rule is basically a tool for reducing error. We (as supervisors of the network) know which mapping function we are training the network to compute. So, we can measure the discrepancy between the output that the network actually produces and the output that it is supposed to produce. We can label that discrepancy  $\delta$  (small delta). It will be a number – the number reached by subtracting the actual output from the correct output. So:

$$\delta = \text{INTENDED OUTPUT} - \text{ACTUAL OUTPUT}$$

Suppose, for example, that we are trying to produce a network that functions as an AND-gate. This means that, when the inputs each have value 1, the desired output is 1 (since A AND B is true in the case where A is true and B is true). If the output that the network actually produces is 0, then  $\delta = 1$ . If, in contrast, the desired output is 0 and the actual output is 1, then  $\delta = -1$ .

It is standard when constructing neural networks to specify a learning rate. This is a constant number between 0 and 1 that determines how large the changes are on each trial.



We can label the learning rate constant  $\varepsilon$  (epsilon). The perceptron convergence rule is a very simple function of  $\delta$  and  $\varepsilon$ .

If we use the symbol  $\Delta$  (big delta) to indicate the adjustment that we will make after each application of the rule, then the perceptron convergence rule can be written like this (remembering that  $T$  is the threshold;  $I_i$  is the  $i$ th input; and  $W_i$  is the weight attached to the  $i$ th input):

$$\begin{aligned}\Delta T &= -\varepsilon \times \delta \\ \Delta W_i &= \varepsilon \times \delta \times I_i\end{aligned}$$

There are multiple adjustments here, one for the threshold and one for each of the weights. Let's look at how the equations work to see how the network learns by making these adjustments.

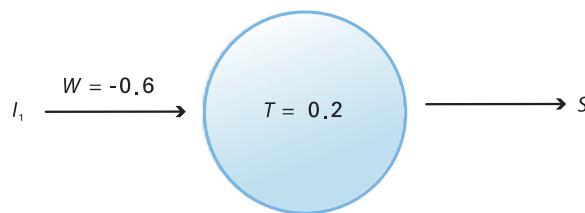
Suppose  $\delta$  is positive. This means that our network has undershot (because it means that the correct output is greater than the actual output). Since the actual output is weaker than required we can make two sorts of changes in order to close the gap between the required output and the actual output. We can *decrease* the threshold and we can *increase* the weights. This is exactly what the perceptron convergence rule tells us to do. We end up decreasing the threshold because when  $\delta$  is positive,  $-\varepsilon \times \delta$  is negative. And we end up increasing the weights, because  $\varepsilon \times \delta \times I_i$  comes out positive when  $\delta$  is positive.



### Exercise 5.5 How does the perceptron convergence rule work when the network overshoots?

Here is an example. Let's consider the very simple single-layer network depicted in Figure 5.7. This network only takes one input and so we only have one weight to worry about. We can take the starting weight to be  $-0.6$  and the threshold to be  $0.2$ . Let's set our learning constant at  $0.5$  and use the perceptron learning rule to train this network to function as a NOT-gate.

Suppose that we input a  $1$  into this network (where, as before,  $1$  represents TRUE and  $0$  represents FALSE). The total input is  $1 \times -0.6 = -0.6$ . This is below the threshold of  $0.2$  and so the output signal is  $0$ . Since this is the desired output we have  $\delta = 0$  and so no learning takes place (since  $\Delta T = -\varepsilon \times \delta = -0.5 \times 0 = 0$ , and  $\Delta W$  also comes out as  $0$ ). But if we input a  $0$  then we get a total input of  $0 \times -0.6 = 0$ . Since this is also below the threshold the output signal is  $0$ . But this is not the desired output, which is  $1$ . So, we can calculate



**Figure 5.7** The starting configuration for a single-layer network being trained to function as a NOT-gate through the perceptron convergence rule. It begins with a weight of  $-0.6$  and a threshold of  $0.2$ .

$\delta = 1 - 0 = 1$ . This gives  $\Delta T = -0.5 \times 1 = -0.5$  and  $\Delta W = 0.5 \times 1 \times 0 = 0$ . This changes the threshold (to  $-0.3$ ) and leaves the weight unchanged.

This single application of the perceptron convergence rule is enough to turn our single-unit network with randomly chosen weight and threshold into a NOT-gate. If we input a 1 into the network then the total input is  $1 \times -0.6 = -0.6$ , which is below the threshold. So the output signal is 0, as required. And if we input a 0 into the network then the total input is  $0 \times -0.6 = 0$ , which is above the threshold of  $-0.3$ . So, the output signal is 1, as required. In both cases we have  $\delta = 0$  and so no further learning takes place. The network has *converged* on a solution.

The perceptron convergence rule is very powerful. In fact, it can be proved (although we shan't do so here) that applying the rule is guaranteed to converge on a solution in every case that a solution exists. But can we say anything about when there is no solution – and hence about which functions a network can learn to compute via the perceptron convergence rule and which will forever remain beyond its reach? It turns out that there is a relatively simple way of classifying the functions that a network can learn to compute by applying the perceptron convergence rule. We will see how to do this next.

## Linear Separability and the Limits of Perceptron Convergence

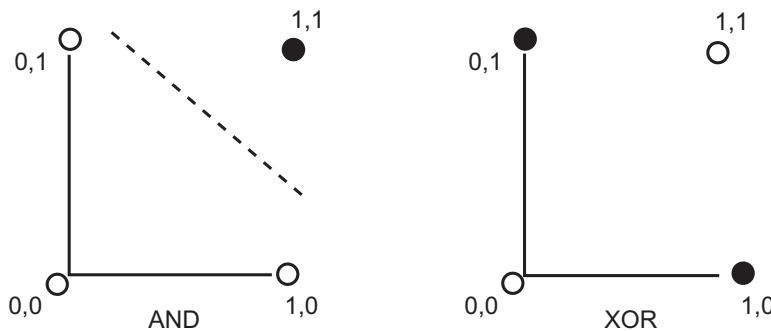
We have seen how our single-layer networks can function as AND-gates, OR-gates, and NOT-gates. And we have also seen an example of how the perceptron convergence rule can be used to train a network with a randomly assigned weight and a randomly assigned threshold to function as a NOT-gate. It turns out that these functions share a common property and that that common property is shared by every function that a single-layer network can be trained to compute. This gives us a very straightforward way of classifying what networks can learn to do via the perceptron convergence rule.

It is easiest to see what this property is if we use a graph to visualize the “space” of possible inputs into one of the gates. Figure 5.8 shows how to do this for two functions. The function on the left is the AND function. On the graph a black dot is used to mark the inputs for which the AND-gate outputs a 1, and a white dot marks the inputs that get a 0. There are four possible inputs and, as expected, only one black dot (corresponding to the case where both inputs have the value TRUE). The graph for AND shows that we can use a straight line to separate out the inputs that receive the value 1 from the inputs that receive the value 0. Functions that have this property are said to be *linearly separable*.



### Exercise 5.6 Draw a graph to show that OR is linearly separable.

Clearly, though, the function on the right is not linearly separable. This is the exclusive-OR function (standardly written as XOR). The OR function that we have been looking at up to now has the value TRUE except when both inputs have the value FALSE. So, A OR B has the value TRUE even when both A and B have the value TRUE. This is not how the word “or” often works in English. If I am offered a choice between A or B it often means that I have to



**Figure 5.8** Graphical representations of the AND and XOR (exclusive-OR) functions, showing the linear separability of AND. Each of the four circles marked on the graph represents a possible combination of input truth values (as fixed by its coordinates). The circle is colored black just if the function outputs 1 at that point.

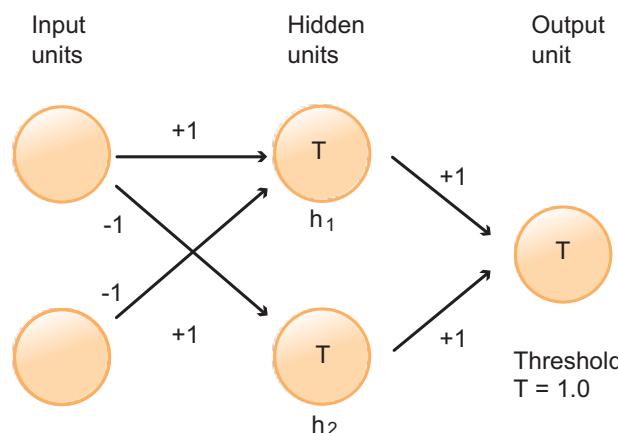
choose one, but not both. This way of thinking about “or” is captured by the function XOR. A XOR B has the value TRUE only when exactly one of A and B has the value TRUE.

No straight line separates the black dots from the white dots in the graph of XOR. This means that XOR is not linearly separable. It turns out, moreover, that XOR cannot be represented by a single-layer network. This is easier to see if we represent XOR in a truth table. The table shows what the output is for each of the four different possible pairs of inputs – as usual, 1 is the TRUE input and 0 is the FALSE input.

$I_1$	$I_2$	OUTPUT
0	0	0
0	1	1
1	0	1
1	1	0

Now, think about how we would need to set the weights and the threshold to get a single-layer network to generate the right outputs. We need the network to output a 1 when the first input is 0 and the second input is 1. This means that  $W_2$  (the weight for the second input) must be such that  $1 \times W_2$  is greater than the threshold. Likewise, for the case where the first input is 1 and the second input is 0. In order to get this to come out right we need  $W_1$  to be such that  $1 \times W_1$  is greater than the threshold. But now, with the weights set like that, it is inevitable that the network will output a 1 when both inputs are 1 – if each input is weighted so that it exceeds the threshold, then it is certain that adding them together will exceed the threshold. In symbols, if  $W_1 > T$  and  $W_2 > T$ , then it is inevitable that  $W_1 + W_2 > T$ .

So, XOR fails to be linearly separable and is also not computable by a single-layer network. In fact, there is a general principle here. The class of Boolean functions that can be computed by a single-unit network is precisely the class of linearly separable functions.



**Figure 5.9** A multilayer network representing the XOR (exclusive-OR) function. Note that, unlike the single-layer perceptrons that we have been considering up to now, this network has three layers. One of these layers is a hidden layer – it receives inputs only indirectly from other units. (Adapted from McLeod, Plunkett, and Rolls 1998)

This was proved by Marvin Minsky and Seymour Papert in a very influential book entitled *Perceptrons* that was published in 1969.

But why does this matter? It is not too hard to construct an artificial network that will compute XOR. Figure 5.9 shows a network that will do the job. It is what is known as a *multilayer network*. Up to now we have been looking at single-layer networks. The units in single-layer networks receive inputs directly. Multilayer networks, in contrast, contain units that only receive inputs indirectly. These are known as *hidden units*. The only inputs they can receive are outputs from other units.



**Exercise 5.7** There are two binary Boolean functions that fail to be linearly separable. The second is the reverse of XOR, which assigns 1 where XOR assigns 0 and 0 where XOR assigns 1. Construct a network that computes this function.

The presence of hidden units is what allows the network in Figure 5.9 to compute the XOR function. The reason a single-layer network cannot compute XOR is that it can only assign one weight to each input. This is why a network that outputs 1 when the first input is 1 and outputs 1 when the second input is 1 has to output 1 when both inputs are 1. This problem goes away when a network has hidden units. Each input now has its own unit and each input unit is connected to two different output units. This means that two different weights can now be assigned to each input.

Multilayered networks can compute any computable function – not just the linearly separable ones. But what stopped researchers in their tracks in 1969 was the fact that they had no idea how to train multilayered networks. The great breakthrough in neural network modeling came with the discovery of an algorithm for training multilayer networks.



**Exercise 5.8** Why can't the perceptual convergence rule be applied to multilayer networks?



## 5.3 Multilayer Networks

Let's start with some basic facts about multilayer networks. Multilayer networks are organized into different layers. Each layer contains a number of units, typically not connected to each other. All networks contain an input layer, an output layer, and a number (possibly 0) of what are called *hidden layers*. The hidden layers are so called because they are connected only to other network units. They are hidden from the "outside world."

Information enters the network via the input layer. Each unit in the input layer receives a certain degree of activation, which we can represent numerically. Each unit in the input layer is connected to each unit in the next layer. Each connection has a weight, again representable numerically. The most common neural networks are *feedforward* networks. As the name suggests, activation spreads forward through the network. There is no spread of activation between units in a given layer, or backward from one layer to the previous layer.

The spread of activation through a multilayer network is illustrated in Figure 5.10, which illustrates a sample hidden unit in a simple network with only one layer of hidden units. (Note that the diagram follows the rather confusing notation standard in the neural network literature.)

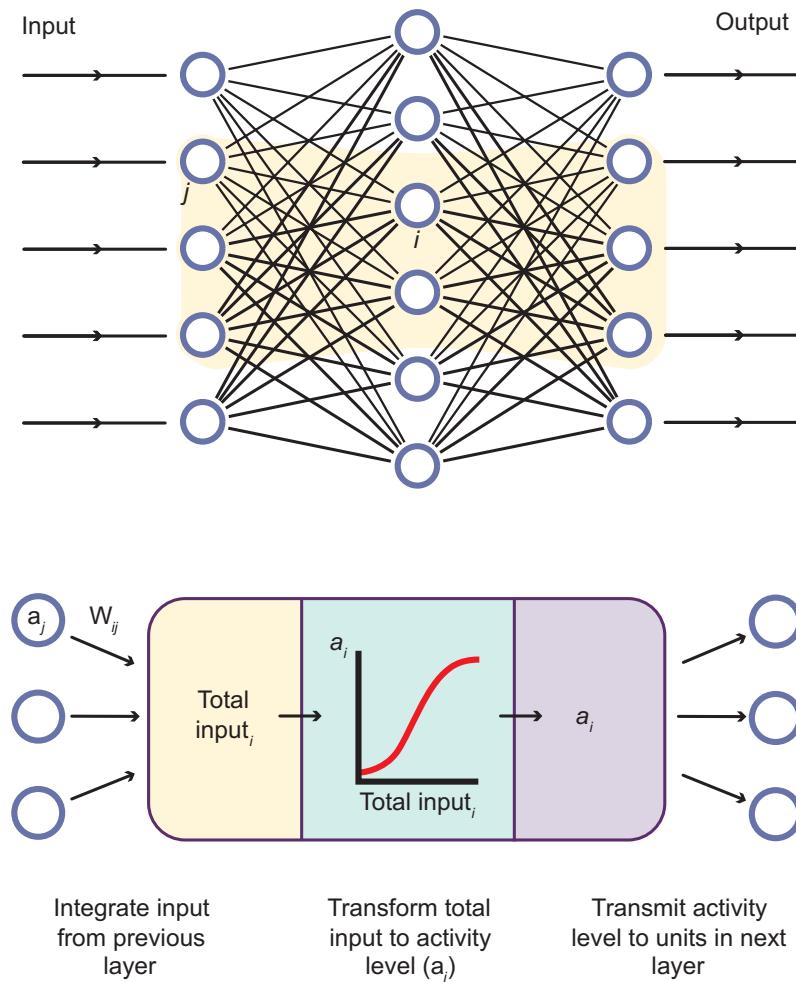
The usual practice is to label a particular unit with the subscript  $i$ . So, we write the name of the unit as  $u_i$ . If we want to talk about an arbitrary unit from an *earlier* layer connected to  $u_i$ , we label that earlier unit with the subscript  $j$  and write the name of the unit as  $u_j$ . Just to make things as difficult as possible, when we label the weight of the connection from  $u_j$  to  $u_i$  we use the subscript  $ij$ , with the label of the later unit coming first. So,  $w_{ij}$  is the weight of the connection that runs *from*  $u_j$  *to*  $u_i$ .

As we see in the figure, our sample unit  $u_i$  integrates the activation it receives from all the units in the earlier layer to which it is connected. Assume that there are  $n$  units connected to  $u_i$ . Multiplying each by the appropriate weight and adding the resulting numbers all together gives the total input to the unit – which we can write as total input ( $i$ ). If we represent the activation of each unit  $u_j$  by  $a_j$ , then we can write down this sum as

$$\text{Total input} = \sum_{j=1}^N w_{ij} a_j$$

We then apply the activation function to the total input. This will determine the unit's activity level, which we can write down as  $a_i$ . In the figure the activation function is a sigmoid function. This means that  $a_i$  is low when total input ( $i$ ) is below the threshold. Once the threshold is reached,  $a_i$  increases more or less proportionally to total input. It then levels out once the unit's ceiling is reached.

Once we understand how a single unit works it is straightforward to see how the whole network functions. We can think of it as a series of  $n$  time steps where  $n$  is the number of layers (including the input, hidden, and output layers). In the first time step every unit in the input layer is activated. We can write this down as an ordered series of numbers – what mathematicians call a *vector*. At step 2 the network calculates the activation level of each



**Figure 5.10** The computational operation performed by a unit in a connectionist model. Upper: General structure of a connectionist network. Lower: A closer look at unit  $i$ . Its operation can be broken into three steps: (1) Integrate all the inputs from the previous layer to create a total input. (2) Use an activation function to convert the total input to an activity level. (3) Output the activity level as input to units in the next layer. (Adapted from McLeod, Plunkett, and Rolls 1998)

unit in the first hidden layer, by the process described in the previous paragraph. This gives another vector. And so on until at step  $n$  the network has calculated the activation level of each unit in the output layer to give the output vector.

## The Backpropagation Algorithm

So that's how activation spreads through a multilayer network. But where do the weights come from? How does the network learn to solve a particular problem, whether it is computing the XOR function or distinguishing between mines and rocks (as in the network we looked at in Section 3.3).



This brings us to the *backpropagation algorithm*. The basic idea is that error is propagated backward through the network from the output units to the hidden units. Recall the basic problem for training multilayer networks. We know what the target activation levels are for the output units. We know, for example, that a network computing XOR should output 0 when the inputs are both 1. And we know that a mine/rock detector should output (1, 0) when its inputs correspond to a mine and (0, 1) when its inputs correspond to a rock. Given this we can calculate the degree of error in a given output unit. But since we don't know what the target activation levels are for the hidden units we have no way of calculating the degree of error in a given hidden unit. And that seems to mean that we have no way of knowing how to adjust the weights of connections to hidden units.

The backpropagation algorithm solves this problem by finding a way of calculating the error in the activation level of a given hidden unit even though there is no explicit activation level for that unit. The basic idea is that each hidden unit connected to an output unit bears a degree of "responsibility" for the error of that output unit. If, for example, the activation level of an output unit is too low, then this can only be because insufficient activation has spread from the hidden units to which it is connected. This gives us a way of assigning error to each hidden unit. In essence, the error level of a hidden unit is a function of the extent to which it contributes to the error of the output unit to which it is connected. Once this degree of responsibility, and consequent error level, is assigned to a hidden unit, it then becomes possible to modify the weights between that unit and the output unit to decrease the error.

This method can be applied to as many levels of hidden units as there are in the network. We begin with the error levels of the output units and then assign error levels to the first layer of hidden units. This allows the network both to modify the weights between the first layer of hidden units and the output units and to assign error levels to the next layer of hidden units. And so the error is *propagated* back down through the network until the input layer is reached. It is very important to remember that activation and error travel through the network in opposite directions. Activation spreads forward through the network (at least in *feed forward* networks), while error is propagated backward.

## How Biologically Plausible Are Neural Networks?

There are certainly some obvious and striking dissimilarities at many different levels between neural networks and the brain. For example –

- Whereas neural network units are all homogeneous, there are many different types of neuron in the brain – twelve different types in the neocortex alone.
- Brains are nowhere near as massively parallel as typical neural networks. Each cortical neuron is connected to a roughly constant number of neurons (approximately 3 percent of the neurons in the surrounding square millimeter of cortex).
- The scale of connectionist networks seems wrong. Each cortical column consists of a population of highly interconnected neurons with similar response properties. A single cortical column cuts vertically across a range of horizontal layers (*laminae*) and can contain

as many as 200,000 neurons – whereas even the most complicated artificial neural networks rarely have more than 5,000 units.

There are even more striking differences when it comes to learning and training:

- Neural networks learn by modifying connection weights and even in relatively simple networks this requires hundreds and thousands of training cycles. (But still – the principal reason why training a network takes so long is that networks tend to start with a random assignment of weights and this is not something one would expect to find in a well-designed brain.)
- There is no evidence that anything like the backpropagation of error takes place in the brain. Researchers have failed to find any neural connections that transmit information about error.
- Most neural networks are supervised networks and require detailed information about the extent of the error at each output unit. But very little biological learning seems to involve this sort of detailed feedback.

It is important to keep these arguments in perspective, however. There are learning algorithms that are more biologically plausible than backpropagation, such as *local algorithms*.

In local learning algorithms an individual unit's weight changes directly as a function of the inputs to and outputs from that unit. Thinking about it in terms of neurons, the information for changing the weight of a synaptic connection is directly available to the presynaptic axon and the postsynaptic dendrite. The Hebbian learning rule that we briefly looked at earlier is an example of a local learning rule. Neural network modelers think of it as much more biologically plausible than the backpropagation rule.

Local learning algorithms, are often used in networks that learn through unsupervised learning. The backpropagation algorithm requires very detailed feedback, as well as a way of spreading an error signal back through the network. *Competitive networks*, in contrast, do not require any feedback at all. There is no fixed target for each output unit and there is no external teacher. What the network does is classify a set of inputs in such a way that each output unit fires in response to a particular set of input patterns.

This works because of inhibitory connections between the output units, in contrast to standard feedforward networks, where there are typically no connections between units in a single layer. These inhibitory connections allow the output units to compete with each other. Each output unit inhibits the other output units in proportion to its firing rate. So, the unit that fires the most will win the competition. Only the winning unit is “rewarded” (by having its weights increased). This increase in weights makes it more likely to win the competition when the input is similar. The end result is that each output ends up firing in response to a set of similar inputs.

Competitive networks are particularly good at classification tasks, which require detecting similarities between different input patterns. They have been used, for example, to model visual pattern recognition. One of the amazing properties of the visual system is its ability to recognize the same object from many different angles and perspectives. There are several competitive network models of this type of *position-invariant object recognition*, including the VisNet model of visual processing developed by Edmund Rolls and



T. T. Milward. VisNet is designed to reproduce the flow of information through the early visual system (as sketched in Section 3.2). It has different layers intended to correspond to the stages from area V1 to the inferior temporal cortex. Each layer is itself a competitive network, learning by a version of the Hebbian rule.

In short, there are many ways of developing the basic insights in neural network models that are more biologically plausible than standard feedforward networks that require detailed feedback and a mechanism for the backpropagation of error. And in any case, neural network models should be judged by the same criteria as other mathematical models. In particular, the results of the network need to mesh reasonably closely with what is known about the large-scale behavior of the cognitive ability being modeled. So, for example, if what is being modeled is the ability to master some linguistic rule (such as the rule governing the formation of the past tense), one would expect a good model to display a learning profile similar to that generally seen in the average language learner. In Chapter 10 we will look at two examples of models that do seem very promising in this regard. First, though, we need to make explicit some of the general features of the neural network approach to information processing.



## 5.4

## Information Processing in Neural Networks: Key Features

Stepping back from details of specific networks and learning rules, all neural networks share some very general characteristics that distinguish them from physical symbol systems.



### Distributed Representations

According to the physical symbol system hypothesis, representations are distinct and identifiable components in a cognitive system. This need not be true in artificial neural networks. There are some networks for which it holds. These are called *localist* networks. What distinguishes localist networks is that each unit codes for a specific feature in the input data. We might think of the individual units as analogs of concepts. They are activated when the input has the feature encoded that the unit encodes. The individual units work as simple feature-detectors. There are many interesting things that can be done with localist networks. But the artificial neural networks that researchers have tended to find most exciting have typically been *distributed* networks rather than localist ones. Certainly, all the networks that we have looked at in this chapter have been distributed.

The information that a distributed network carries is not located in any specific place. Or rather, it is distributed across many specific places. A network stores information in its pattern of weights. It is the particular pattern of weights in the network that determines what output it produces in response to particular inputs. A network learns by adjusting its weights until it settles into a particular configuration – hopefully the configuration that produces the right output! The upshot of the learning algorithm is that the network's "knowledge" is distributed across the relative strengths of the connections between different units.

## No Clear Distinction between Information Storage and Information Processing

According to the physical symbol system hypothesis all information processing is rule-governed symbol manipulation. If information is carried by symbolic formulas in the language of thought, for example, then information processing is a matter of transforming those formulas by rules that operate only on the formal features of the formulas. In the last analysis, information is carried by physical structures and the rules are rules for manipulating those symbol structures. This all depends upon the idea that we can distinguish within a cognitive system between the representations on which the rules operate and the rules themselves – just as, within a logical system such as the propositional or predicate calculus, we can distinguish between symbolic formulas and the rules that we use to build those symbolic formulas up into more complex formulas and to transform them.



**Exercise 5.9** Look back at Box 4.1 and Figure 4.2 and explain how and why the distinction between rules and representations is central to the physical symbol system and language of thought hypotheses.

Consider how AND might be computed according to the physical symbol system hypothesis. A system for computing AND might take as its basic alphabet the symbol 0 and the symbol 1. The inputs to the system would be pairs of symbols and the system would have built into it rules to ensure that when the input is a pair of 1s, the system outputs a 1, while in all other cases, it outputs a 0. What might such a rule look like?

Well, we might think about the system along the lines of a Turing machine (as illustrated in Section 1.2). In this case the inputs would be symbols written on two squares of a tape. Assume that the head starts just to the left of the input squares. The following program will work.

*Step 1* Move one square R.

*Step 2* If square contains “1,” then delete it, move one square R and go to Step 6.

*Step 3* If square contains “0,” then delete it, move one square R and go to Step 4.

*Step 4* Delete what is in square and write “0.”

*Step 5* Stop.

*Step 6* If square contains “0,” then stop.

*Step 7* If square contains “1,” then stop.

The tape ends up with a 1 on it only when the tape started out with two 1s on it. If the tape starts out with one or more 0s on it then it will stop with a 0. The final state of the tape is



reached by transforming the initial symbol structure by formal rules, exactly as required by the physical symbol system hypothesis. And the rules are completely distinct from the symbols on which they operate.



### Exercise 5.10 Write a program that will compute the function XOR.

There is no comparable distinction between rules and representations in artificial neural networks. The only rules are those governing the spread of activation values forward through the network and those governing how weights adjust. Look again at the network computing XOR and think about how it works. If we input two 1s into the network (corresponding to a pair of propositions, both of which are true), then the information processing in the network proceeds in two basic stages. In the first-stage activation spreads from the input layer to the hidden layer and both hidden units fire. In the second stage, activation spreads from the hidden units to the output unit and the output unit fires.

The only rules that are exploited are, first, the rule for calculating the total input to a unit and, second, the rule that determines whether a unit will fire for a given total input (i.e., the activation function). But these are exactly the same rules that would be activated if the network were computing AND or OR. These “updating rules” apply to all feedforward networks of this type. What distinguishes the networks are their different patterns of weights. But a pattern of weights is not a rule, or an algorithm of any kind. Rather a particular pattern of weights is what results from the application of one rule (the learning algorithm). And it is one of the inputs into another rule (the updating algorithm).



## The Ability to Learn from “Experience”

Of course, talk of neural networks learning from experience should not be taken too seriously. Neural networks do not experience anything. They just receive different types of input. But the important point is that they are not fixed in how they respond to inputs. This is because they can change their weights. We have looked at several different ways in which this can take place – at several different forms of learning algorithm. Supervised learning algorithms, such as the backpropagation algorithm, change the weights in direct response to explicit feedback about how the network’s actual output diverges from intended output. But networks can also engage in unsupervised learning (as we saw when we looked briefly at competitive networks). Here the network imposes its own order on the inputs it receives, typically by means of a local learning algorithm, such as some form of Hebbian learning.

This capacity to learn makes neural networks a powerful tool for modeling cognitive abilities that develop and evolve over time. We will look at examples of how this can be done later on, particularly in Chapters 10 and 12.



## Summary

This chapter has explored a way of thinking about information processing very different from the physical symbol system hypothesis discussed in Chapter 4. Connectionist neural networks are constructed from individual units that function as highly idealized neurons. We looked at two very different types of network. In the first part of the chapter we looked at single-layer networks and saw how they can learn via the perceptron convergence rule. Unfortunately, single-layer networks are limited in the functions that they can compute. It has been known for a long time that multilayer networks built up from single-layer networks can compute any function that can be computed by a digital computer, but it was not until the emergence of the backpropagation learning algorithm that it became possible to train multilayer neural networks. The chapter ended by considering the biological plausibility of neural networks and summarizing some of the crucial differences between artificial neural networks and physical symbol systems.

## Checklist

### Neurally Inspired Information Processing

- (1) A fundamental question in thinking about how the brain processes information is how the activities of large populations of neurons give rise to complex sensory and cognitive abilities.
- (2) Existing techniques for directly studying the brain do not allow us to study what happens inside populations of neurons.
- (3) Computational neuroscientists use mathematical models (neural networks) to study populations of neurons.
- (4) These neural networks are made up of units loosely based on biological neurons. Each unit is connected to other units so that activation levels can be transmitted between them as a function of the strength of the connection.

### Single-Layer Networks

- (1) We can use single-layer networks to compute some Boolean functions, in particular AND, OR, and NOT.
- (2) Any digital computer can be simulated by a network of single-layer networks appropriately chained together.
- (3) Single-layer networks can learn by adjusting their weights to minimize their degree of error (the  $\delta$  signal) according to the perceptron convergence rule.
- (4) Single-layer networks can only learn to compute functions that are linearly separable.



### Multilayer Networks

- (1) Multilayer networks have hidden units that are neither input units nor output units.
- (2) The presence of hidden units enables multilayer networks to learn to compute functions that cannot be learned by single-layer networks (including functions that are not linearly separable).
- (3) The backpropagation learning algorithm for multilayer networks adjusts the weights of hidden units as a function of how “responsible” they are for the error at the output units.

### Biological Plausibility

- (1) Neural network units are much more homogeneous than real neurons. And real neural networks are likely to be both much larger and less parallel than network models.
- (2) The backpropagation algorithm is not very biologically plausible. There is no evidence that error is propagated backward in the brain. And nature rarely provides feedback as detailed as the algorithm requires.
- (3) However, there are other learning algorithms. Competitive networks using Hebbian learning do not require explicit feedback, and there is evidence for local learning in the brain.

### Information Processing in Neural Networks

- (1) Representation in neural networks is distributed across the units and weights, rather than being encoded in discrete symbol structures, as in the physical symbol system hypothesis.
- (2) There are no clear distinctions to be drawn within neural networks either between information storage and information processing or between rules and representations.
- (3) Neural networks are capable of sophisticated forms of learning, which makes them particularly suitable for modeling how cognitive abilities are acquired and how they evolve.

## Further Reading

The *Handbook of Brain Theory and Neural Networks* (Arbib 2003) is the most comprehensive single-volume source for different types of computational neuroscience and neural computing, together with entries on neuroanatomy and many other “neural topics.” It contains useful introductory material and “road maps.” Stein and Stoodley 2006 and Trappenberg 2010 are user-friendly introductions to neuroscience and computational neuroscience, respectively. Arbib 1987 surveys the theoretical issues in modeling the brain from a mathematical perspective. Glass 2016 is a neuroscience-inspired cognitive psychology textbook.

The classic sources for connectionism are the two volumes of Rumelhart, McClelland, and the PDP Research Group 1986. Churchland and Sejnowski 1992 is an early manifesto for computational neuroscience. See also Bechtel and Abrahamsen 2002 and the relevant chapters of Dawson 1998. There are useful article-length presentations in Rumelhart 1989 (reprinted in

Haugeland 1997) and Churchland 1990b (reprinted in Cummins and Cummins 2000). McLeod, Plunkett, and Rolls 1998 covers both the theory of neural networks and their modeling applications, including the VisNet model of visual processing originally presented in Rolls and Milward 2000. The first chapter is reprinted in Bermúdez 2006. Dawson 2005 is a "hands-on" introduction to connectionist modeling. For a survey of applications of connectionist networks in cognitive psychology, see Houghton 2005. See also Thomas and McClelland's chapter on connectionist modeling in Sun 2008. A more recent discussion of connectionism can be found in McClelland et al. 2010, with commentary and target articles from others in the same issue. Connectionism went out of fashion for a few years but has recently seen a resurgence in the context of deep machine learning. See the Further Reading section in Chapter 12 for references.

The biological plausibility of artificial neural networks has been much discussed, and researchers have developed a number of learning algorithms that are less biologically implausible than the backpropagation algorithm. O'Reilly and Munakata 2000 is a good place to start in finding out about these. Warwick 2012 is a more recent alternative. See Bowers 2009 and Plaut and McClelland 2010 for an exchange concerning biological plausibility as well as local and distributed representations. The perceptron convergence learning rule discussed in Section 8.2 is also known as the delta rule. It is very closely related to the model of associative learning in classical (Pavlovian) conditioning independently developed by the psychologists Robert Rescorla and Allen Wagner in the 1970s. For more on reward learning and the delta rule, see chapter 6 of Trappenberg 2010. The *Encyclopedia of Cognitive Science* also has an entry on perceptrons (Nadel 2005). For more on McCullough and Pitts, see chapter 2 of Arbib 1987 and Piccinini 2004, as well as Schlatter and Aizawa 2008.

One of the key distinguishing features of neural networks is that their "knowledge" is distributed across units and weights. This raises a number of issues, both practical and theoretical. Rogers and McClelland 2004 develops a distributed model of semantic knowledge. Philosophers have explored the relation between distributed representations and standard ways of thinking about propositional attitudes and mental causation. Some of the points of contact are explored in Clark 1989, 1993. Macdonald and Macdonald 1995 collects some key papers, including an important debate between Smolensky and Fodor about the structure of connectionist networks. Other collections include Davis 1993 and Ramsey, Stich, and Rumelhart 1991.

Not all neural networks are distributed. There are also localist networks. Whereas in distributed networks, it is typically not possible to say what job an individual unit is doing (and when it is possible, it usually requires knowing a lot about what other units are doing), units in localist networks can be interpreted independently of the states of other units. For a robust defense of the localist approach, see Page 2000 and the papers in Grainger and Jacobs 1998.

One topic not discussed in the text is the computational power of artificial neural networks. It is sometimes suggested that connectionist networks are computationally equivalent to digital computers (in virtue of being able to compute all Turing-computable functions), which might be



taken to indicate that connectionist networks are simply implementations of digital computers. The implementation thesis is canvassed both by opponents of connectionism (Fodor and Pylyshyn 1988) and by leading connectionist modelers (Hinton, McClelland, and Rumelhart 1986). Siegelmann and Sontag 1991 present a neural network that can simulate a universal Turing machine. For skeptical discussion, see Hadley 2000.





## CHAPTER SIX

# Applying Dynamical Systems Theory to Model the Mind

### OVERVIEW 149

#### 6.1 Cognitive Science and Dynamical Systems 149

What Are Dynamical Systems? 150  
The Dynamical Systems Hypothesis: Cognitive Science without Representations? 153

#### 6.2 Applying Dynamical Systems: Two Examples from Child Development 158

Two Ways of Thinking about Motor Control 159  
Dynamical Systems and the A-Not-B Error 161  
Assessing the Dynamical Systems Approach 166



## Overview

We have been exploring the basic idea that cognition is information processing. We have looked at different ways of thinking about information processing – the physical symbol system hypothesis and the neural networks model. These two approaches are both committed to thinking of cognition as essentially a process of transforming representational states that carry information about the agent and about the environment, although they think about these representational states in very different ways.

This chapter introduces a very different way of modeling cognitive abilities. First, we look at how some cognitive scientists have proposed using the mathematical and conceptual tools of dynamical systems theory to model cognitive skills and abilities. As we'll see, dynamical systems models differ in certain fundamental respects from the information-processing models we have been looking at. Then in Section 6.2 we explore two examples of how dynamical systems models can shed light on child development.



6.1

## Cognitive Science and Dynamical Systems

The dynamical systems hypothesis calls for cognitive science to be freed from its dependence on ideas of representation and computation. Its fundamental idea is that we can

understand how organisms respond to the environment and orient themselves in it without assuming that there are internal cognitive systems that carry out specific information-processing tasks. The basic currency of cognitive science is not the information-carrying representation, and nor are computations and algorithms the best way to think about how cognition unfolds.

Instead, the proposal is that cognitive scientists need to use the tools of dynamical systems theory in order to understand how perceivers and agents are embedded in their environments. Dynamical systems have been studied in physics and other natural sciences for many centuries. What's new is the idea that we can understand how cognition works by thinking of cognitive agents as dynamical systems.



## What Are Dynamical Systems?

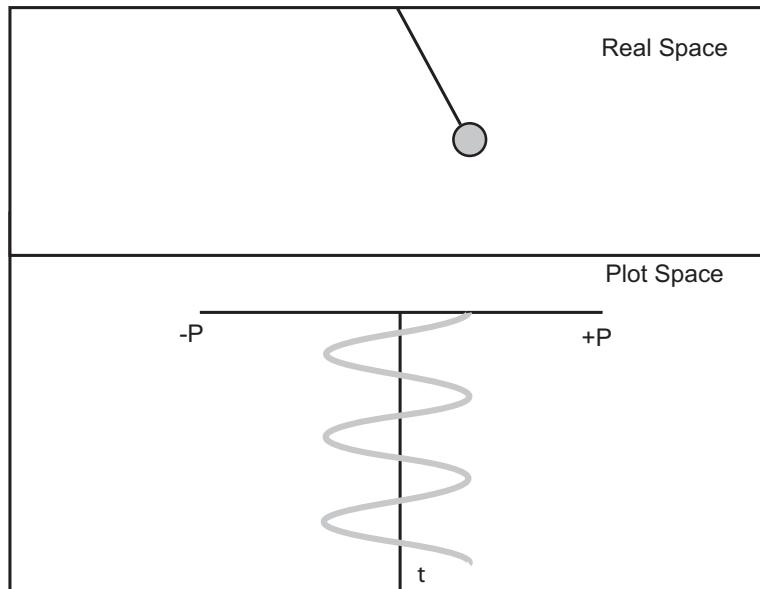
In the broadest sense, a dynamical system is any system that evolves over time in a law-governed way. The solar system is a dynamical system. So are you and I. So is a dripping tap. And so, for that matter, are Turing machines and artificial neural networks. What marks out the dynamical systems hypothesis is the idea that cognitive systems should be studied with the tools of dynamical modeling.

Dynamical modeling exploits powerful mathematical machinery to understand how certain types of natural phenomena evolve over time. Newtonian mechanics is perhaps the most famous example of dynamical modeling, but all dynamical models have certain basic features.

Dynamical models typically track the evolving relationship between a relatively small number of quantities that change over time. They do this using calculus and differential or difference equations. Difference equations allow us to model the evolution of a system that changes in discrete steps. So, for example, we might use difference equations to model how the size of a biological population changes over time – each step being a year, for example. Differential equations, in contrast, allow us to model quantities that change continuously, such as the acceleration of a falling object.

One of the basic theoretical ideas in dynamical systems modeling is the idea of a *state space*. The state space of a dynamical system is a geometric way of thinking about all the possible states that a system can be in. A state space has as many different dimensions as it has quantities that vary independently of each other – as many different dimensions as there are degrees of freedom in the system, in other words. Any state of the dynamical system will involve the system having a particular value in each dimension. And so, we can uniquely identify the state of the system in terms of a particular set of coordinates in the system's state space.

You can think about the state space of an idealized swinging pendulum, for example, as having two dimensions – one corresponding to its angle of displacement from the vertical and one corresponding to its angular velocity. So, every possible state that the pendulum can be in can be represented by a pair of numbers, which in turn correspond to a point in a two-dimensional space.



**Figure 6.1** The trajectory through state space of an idealized swinging pendulum. The pendulum's position is its displacement from the vertical (positive to the right and negative to the left). The time axis goes vertically downward.

If we add another dimension to the state space to represent time then we can start thinking about the evolution of the pendulum in terms of a *trajectory* through state space. A trajectory through state space is simply a sequence of points in the multidimensional space. This sequence of points represents the successive states of the pendulum.

One of the basic aims of dynamical systems modeling is to write equations governing the evolution of the system – that is, governing the different possible trajectories that the system can take through state space, depending upon where the system starts (the system's *initial conditions*).

Let's go back to our idealized simple pendulum, a suspended weight swinging from side to side in an environment with no friction. This is a dynamical system. Its initial condition is its position at the moment it is released and allowed to start swinging. The state space is the different possible positions that the weight can take at a given time. At any given moment, the position of the weight is fixed solely by the pendulum's amplitude (its angle of displacement from its equilibrium position, which is hanging straight down) and the length of time it has been swinging.

So, we can represent the swinging of the pendulum as a trajectory through state space over time. This is represented diagrammatically in Figure 6.1. The state space is two-dimensional, because there are two dimensions of variation. The first dimension of variation is its angular displacement from the vertical equilibrium position. And the second dimension is time, which is represented vertically. So, each state of this simple dynamical system is a position-time pair – a point in the two-dimensional state space.

If we remove some of the simplifying assumptions (by taking friction into account, for example), then our dynamical system becomes more complicated. Since what friction does is decrease velocity through energy loss, we now need to add a third dimension to the state space. This third dimension represents velocity. So now, each state of the system is a position-time-velocity triple – a point in the three-dimensional state space. (See Figure 6.2.)

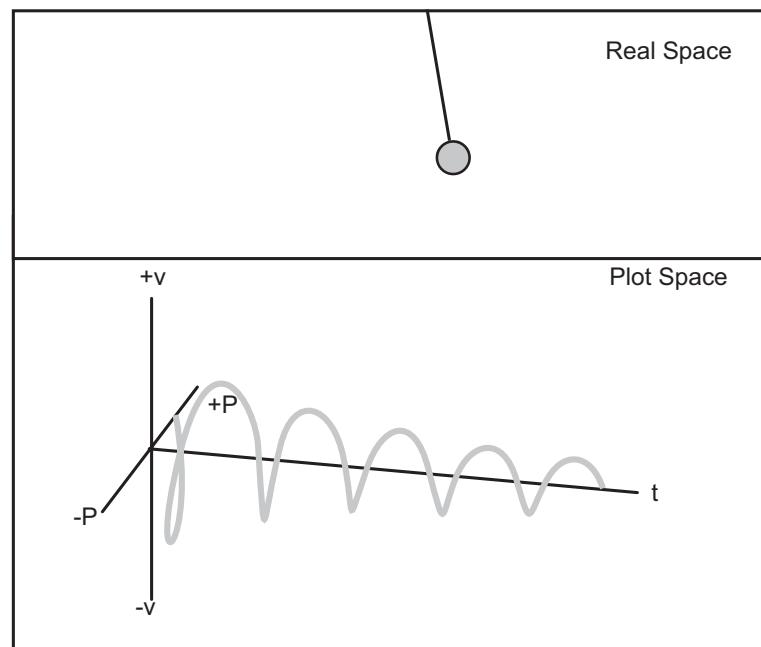


### Exercise 6.1 Explain in your own words the key differences between Figures 6.1 and 6.2.

These two examples illustrate how the evolution of a physical system can be viewed as a trajectory through a multidimensional state space. However, neither is sufficiently detailed to allow us to model the evolution of the system.

In order to model the evolution of the system we need to be able to write equations that describe how a value on one dimension (say, its angle of displacement from the vertical) is determined by values on the other dimensions (say, velocity and time). There is no equation that will fix the position of a swinging pendulum as a function simply of velocity and time (even if we remove the effects of friction). More information is required.

In order to be able to work out where the pendulum will be at a particular moment, we need to take into account two different forces. We need to factor in both the force with which the pendulum is moving and the force of gravity (which acts as a restoring force, counteracting the initial force). When the restoring force exceeds the initial force, then the pendulum starts moving back to equilibrium. And so, it continues to oscillate indefinitely



**Figure 6.2** The state space of a swinging pendulum in a three-dimensional state space. This state space includes a dimension representing velocity, in order to capture energy loss and decreased velocity due to friction.



(in a system where there is no friction). And in order to measure the effects of those forces, you also need to take into account the length of the pendulum.

If we have all this information, then we can start to think about how to write an equation for determining the angle of displacement as a function of these other quantities. It turns out that, with some important simplifying assumptions, the following equation will work when the maximum angle of displacement is relatively small (and there is no friction):

$$A = A_{\text{MAX}} \sin \sqrt{(g/l)} t$$

where  $A$  is the angular displacement;  $A_{\text{MAX}}$  is the maximum angular displacement;  $g$  is the gravitational acceleration; and  $l$  is the length of the pendulum.

The equations get much more complicated when friction is brought back into the picture, and the maximum angle of displacement gets larger. But this should be enough to illustrate the basic idea of a dynamical system and the evolution of a dynamical system can be modeled as a trajectory through state space.

Still, you may reasonably ask, what has this got to do with cognitive science?

## The Dynamical Systems Hypothesis: Cognitive Science without Representations?

To see the relevance of dynamical systems to cognitive science, consider a famous illustration introduced by the philosopher Tim Van Gelder, one of the early proponents of the dynamical systems hypothesis. Van Gelder introduces us to two ways of thinking about an engineering problem whose solution was a vital step in the Industrial Revolution. One way of solving the problem is structurally very similar to the information-processing approach to thinking about how the mind solves problems. The other way, which is how the problem was actually solved, reveals the power of the dynamical systems approach.

For Van Gelder, cognitive scientists are essentially trying to reverse engineer the mind – they are trying to work out how the mind is configured to solve the problems that it deals with. Cognitive scientists have tended to tackle this *reverse engineering* problem in a particular way – by assuming that the mind is an information-processing machine. But what Van Gelder tries to show is that this approach is neither the only way nor the best way. He does this by looking at an example from engineering itself – the Watt governor.

The development of the steam engine is very closely associated with the name of the Scottish engineer James Watt. The first steam engines were only capable of a reciprocating pumping motion. But Watt designed a gearing system to allow steam engines to drive a flywheel and hence to produce rotational power. This gearing system made it possible to use steam engines for weaving, grinding, and other industrial applications.

Unfortunately, there was still a problem. The type of applications for which steam power was needed required the power source to be as uniform as possible. This, in turn, required

the speed of the flywheel to be as constant as possible. But this was very hard to achieve because the speed of the flywheel depended upon two things that were constantly changing – the pressure of the steam driving the engine and the amount of work that the engine was doing. What was needed (and what Watt ended up inventing) was a *governor* that would regulate the speed of the flywheel.

The problem is clear, but how could it be solved? Van Gelder identifies one possible approach. This approach employs the sort of task analysis that is typical of traditional cognitive science, and that is often presented in a boxes-and-arrows diagram. It breaks the task of regulating the speed of the flywheel into a series of subtasks, assumes that each of those subtasks is carried out in separate stages, and works out an algorithm for solving the problem by successively performing the subtasks. This approach gives what Van Gelder terms the *computational governor*, following something like the following algorithm:

- 1 Measure the speed of the flywheel
- 2 Compare the actual speed  $S_1$  against the desired speed  $S_2$
- 3 If  $S_1 = S_2$ , return to step 1
- 4 If  $S_1 \neq S_2$  then
  - (a) measure the current steam pressure
  - (b) calculate the required alteration in steam pressure
  - (c) calculate the throttle adjustment that will achieve that alteration
- 5 Make the throttle adjustment
- 6 Return to step 1

The computational governor has certain features that should be very familiar by now. It is:

*Representational*. It cannot work without some way of representing the speed of the flywheel, the pressure of the steam, and the state of the throttle valve.

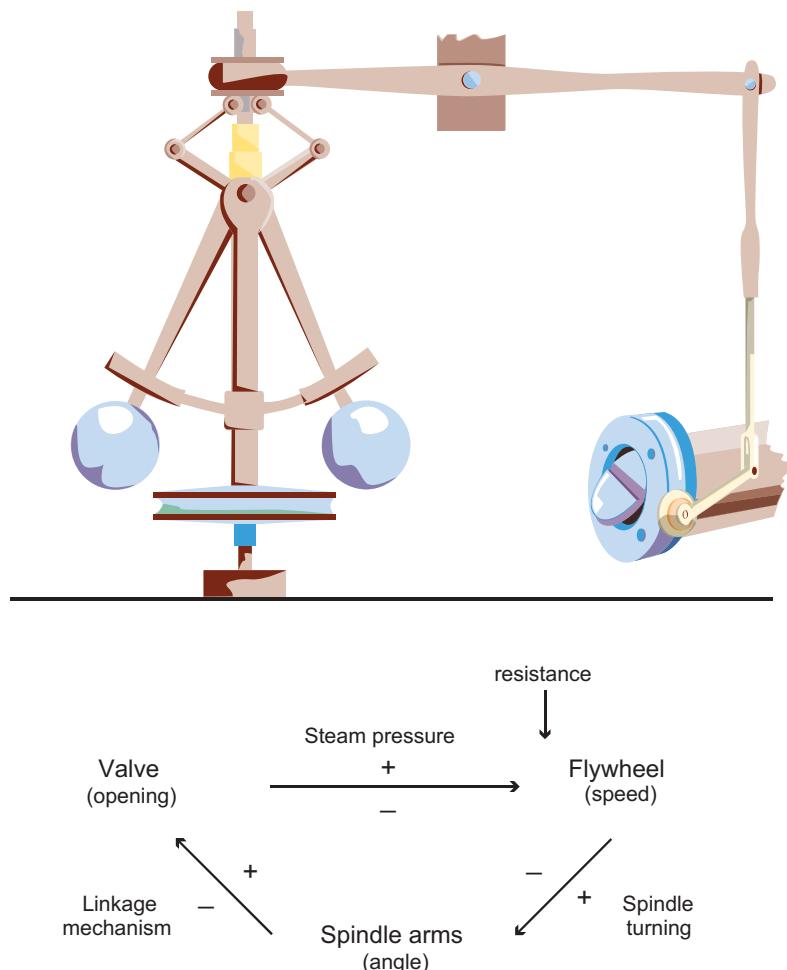
*Computational*. The algorithm is essentially a process for comparing, transforming, and manipulating representations of speed, steam pressure, and so on.

*Sequential*. It works in a discrete, step-by-step manner.

*Decomposable* (or, as Van Gelder puts it, *homuncular*). We can think of it as made up of distinct and semi-autonomous sub-systems, each responsible for a particular sub-task – the speed measurement system, the steam measurement system, the throttle adjustment system, and so on.

So, the computational governor is an application of some of the basic principles that cognitive scientists use to understand how the mind works. The basic fact of the matter, though, is that Watt went about things in a very different way.

The Watt governor, which Watt developed using basic principles already exploited in windmills, has none of the key features of the computational governor. It does not involve



**Figure 6.3** Illustration of the Watt governor, together with a schematic representation of how it works.

representations and hence, as a consequence, cannot be computational. It is not sequential. And it is not decomposable. It is in fact, as Van Gelder points out, a dynamical system that is best studied using the tools of dynamical systems modeling.

The Watt governor is illustrated at the top of Figure 6.3. The flywheel is right at the bottom. Coming up from the flywheel is a rotating spindle. The spindle rotates at a speed determined by the speed of the flywheel. It has two metal balls attached to it. As the speed of the spindle's rotation increases, centrifugal force drives the metal balls upward. As the speed decreases the balls drop down. Watt's key idea was to connect the arms from which the metal balls are suspended directly to the throttle valve for the steam engine. Raising the arms closes down the throttle valve, while the valve is opened up when the arms fall.

This ingenious arrangement allows the governor to regulate the speed by compensating almost instantaneously whether the speed of the flywheel is overshooting or undershooting. The lower part of Figure 6.3 illustrates the feedback loop.

Van Gelder stresses four very important features of the Watt governor:

*Dynamical system:* The best way to understand the Watt governor is through the tools of dynamical systems theory. It is relatively straightforward to write a differential equation that will specify how the arm angle changes as a function of the engine speed. The system is a typical dynamical system because these equations have a small number of variables.

*Time-sensitivity:* The Watt governor is all about timing. It works because fluctuations in the speed of the flywheel are almost instantly followed by variation in the arm angle. The differential equations governing the evolution of the system track the relation over time between flywheel speed and arm angle.

*Coupling:* The Watt governor works because the arm angle, the throttle valve, and the speed of the flywheel are all interdependent. The arm angle is a parameter fixing the speed of the flywheel. But the speed of the flywheel is equally a parameter fixing the angle of the arm. The system as a whole is what dynamical systems theorists call a coupled system characterized by feedback loops.

*Attractor dynamics:* For any given engine speed there is an equilibrium arm angle – an angle that will allow the engine to continue at that speed. We can think about this equilibrium arm angle as an attractor – a point in state space to which many different trajectories will converge. (See Box 6.1.)

So, the Watt governor can be characterized using the tools of dynamical systems theory. It is a coupled system that displays a simple version of attractor dynamics, because it contains basins of attraction (as described in Box 6.1). Unlike the computational governor, it does not involve any representation, computation, or decomposable subsystems. Finally, the Watt governor works in real time. The adjustments are made almost instantaneously, exactly as required. It is very hard to see how the computational governor would achieve this.



### **Exercise 6.2** Explain in your own words the principal differences between the computational governor and the Watt governor.

But again, what has this got to do with the mind? It is not news, after all, that steam engines are not cognitive systems.

Van Gelder and other supporters of the dynamical systems hypothesis argue that the same basic tools that explain how the Watt governor works can be used to illuminate the workings of the mind. But the issue is not just about explanation. Dynamical systems theorists think that the explanations work because they track the basic design principles of the mind. They think not only that the mind is a dynamical system, but also that when we look at the relation between the organism and the environment what we see is a coupled system. The organism–environment complex is a system whose behavior evolves as a function of a small number of variables.



Certainly, the real test of this idea must come in concrete applications. The plausibility of the dynamical systems hypothesis cannot rest solely on an analogy between the mind and a steam engine – however suggestive that analogy may be. Some very exciting work has been done by cognitive scientists on giving dynamical systems models of particular cognitive abilities. Much of the most interesting research has been done on motor skills and

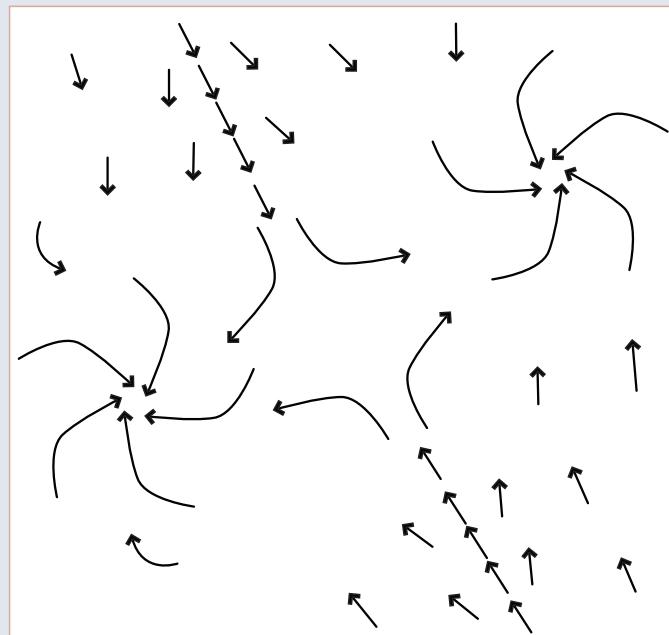
### BOX 6.1 Basins of Attraction in State Space

A particular dynamical system evolves through time along a trajectory in state space. The particular trajectory that it takes is a function of its initial conditions. So, the trajectory of the swinging pendulum, for example, is typically determined by its initial amplitude, together with the way that allowances for friction are built into the system.

But not all regions of state space are created equal. There are some regions of state space to which many different trajectories converge. These are called *basins of attraction*. In the case of a swinging pendulum subject to friction, there is a region of state space to which all trajectories converge – this is the point at which the pendulum is stationary, its equilibrium point.

Many dynamical systems have a number of basins of attraction – these are the *nonlinear dynamical systems*. There is a two-dimensional example in Figure B6.1.

The figure illustrates a range of different possible trajectories. The trajectories are marked by arrows, with the length of the arrow indicating the speed (and hence the strength of the attraction). The state space has two basins of attraction.



**Figure B6.1**

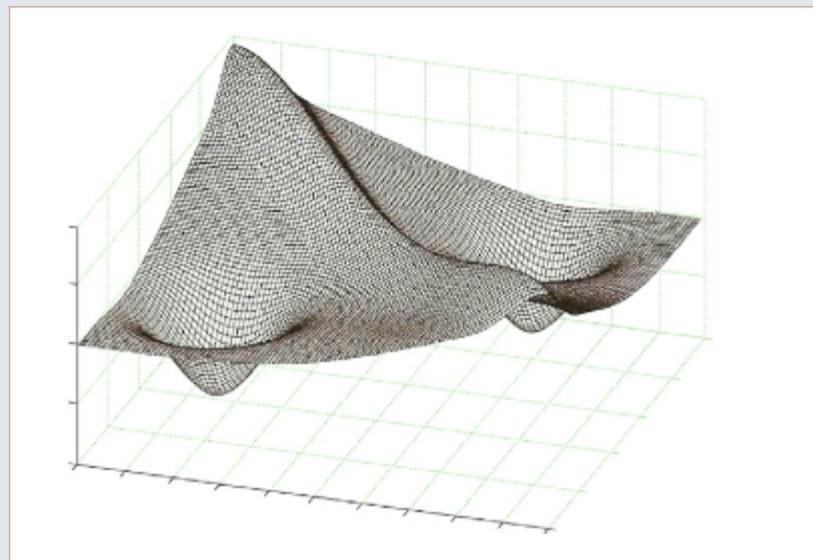
**BOX 6.1 (cont.)****Figure B6.2**

Figure B6.2 gives a different way of representing basins of attraction, in terms of what is often called an energy landscape. This gives a different way of visualizing how a system evolves through state space.

The undulating surface represents the space of possible trajectories. The two basins of attraction are represented by depressions in the surface. Since dynamical systems evolve toward a reduction in energy, trajectories will typically “roll” downhill until they end up in one of the two basins of attraction. In fact, in this particular dynamical system any trajectory must begin on one side of the dividing “ridge” or the other – and so will end up in the corresponding basin of attraction.

motor learning. Dynamical systems theory has proved a powerful tool for understanding how children learn to walk, for example. The next section looks at two applications of the dynamical systems approach to child development.

**6.2**

## Applying Dynamical Systems: Two Examples from Child Development

Dynamical models are extremely time-sensitive, able to track the evolution of a system over time in very fine detail. This suggests that one profitable area to apply them is in



studying how children learn new skills and abilities. In this section we look at two concrete examples of how episodes in child development can be modeled by dynamical systems theory.

## Two Ways of Thinking about Motor Control

Our first example comes from motor control. It has to do with how infants learn to walk. There is a direct analogy with the example of the Watt governor. The dominant approach to understanding how movements are planned and executed is the computational model of motor control. This is the motor control equivalent of the computational governor. The dynamical systems approach offers an alternative – a noncomputational way of thinking about how movements are organized and how motor skills emerge.

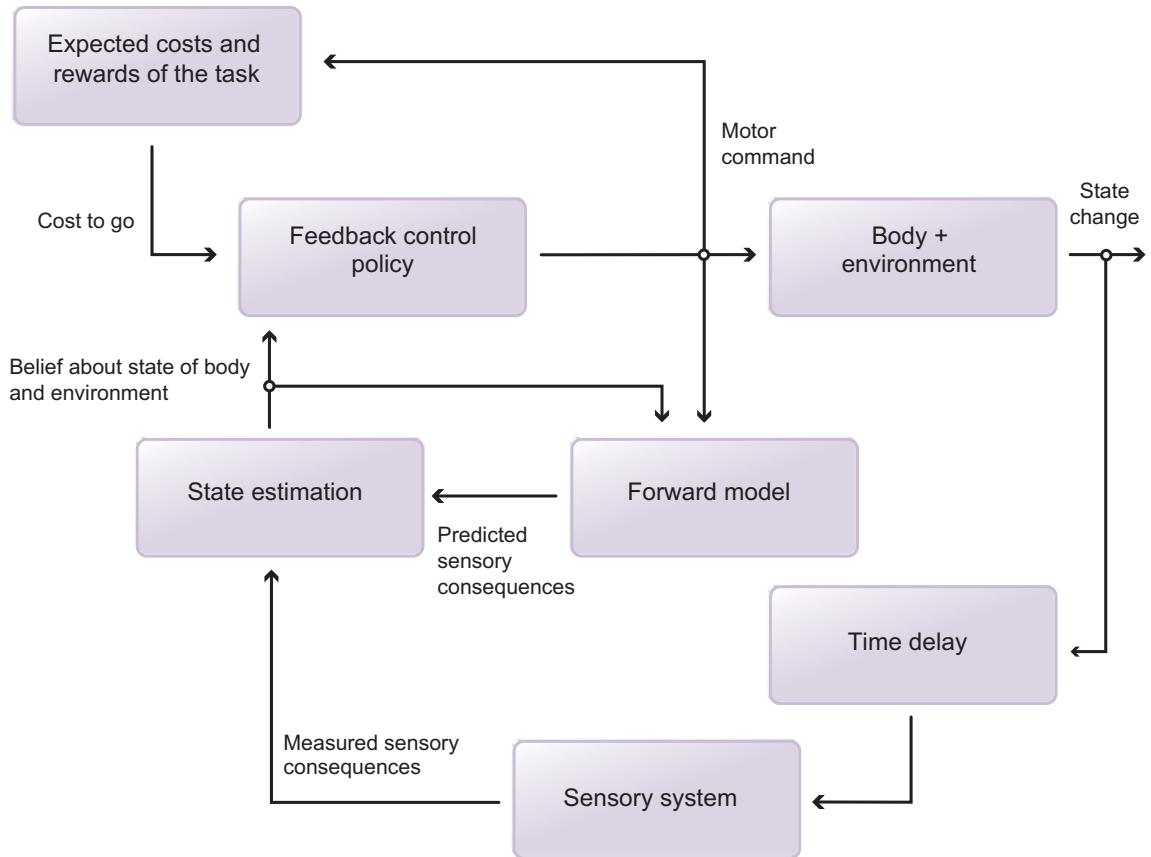
To illustrate the computational model of motor control, consider the movement of reaching for an object. According to the computational model, planning this movement has to begin with the central nervous system calculating both the position of the target object and the position of the hand. These calculations will involve input both from vision and from different types of proprioception (such as sensors in the arm detecting muscle flexion). Planning the movement requires calculating a trajectory from the starting position to the goal position. It also involves computing a sequence of muscle movements that will take the hand along that trajectory. Finally, executing the movement requires calculating changes in the muscle movements to accommodate visual and proprioceptive feedback.

This gives a multistage sequence of computations that seems tailor-made for algorithmic information processing. Figure 6.4 illustrates a computational model of motor control that fits this general description. It is a standard information-processing diagram – which is often called a boxological diagram (because it is drawn in terms of boxes and arrows, with different cognitive tasks assigned to different boxes, and the arrows indicating the flow of information processing).

But the psychologists Esther Thelen and Linda Smith have argued that walking is not a planned activity in the way that many cognitive scientists have assumed, following the computational approach to motor control. It does not involve a specific set of motor commands that “program” the limbs to behave in certain ways. Rather, the activity of walking emerges out of complex interactions between muscles, limbs, and different features of the environment. There are many feedback loops controlling limb movements as a function of variation in both body and environment.

Concrete evidence for Thelen and Smith’s position comes from studies on how infants learn to walk. Most normal infants start learning to walk toward the end of their first year – at around 11 months. For the first few months infants are capable of making stepping movements. They stop making these movements during the so-called nonstepping window. The movements obviously reappear when the infant starts walking.

The traditional explanation for this U-shaped developmental trajectory is that the infant’s initial stepping movements are purely reflexive. They disappear during the



**Figure 6.4** An example of the computational approach to motor control. This model incorporates both forward mechanisms (i.e., mechanisms that make predictions about the sensory consequences of particular movements) and comparator mechanisms (i.e., mechanisms that compare the predictions with actual sensory feedback). (Adapted from Shadmehr and Krakauer 2008)

nonstepping window because the cortex is maturing enough to inhibit reflex responses – but is not sufficiently mature to bring stepping movements under voluntary control.

Thelen and Smith came up with a range of experimental evidence challenging this approach. They discovered that stepping movements could be artificially induced in infants by manipulating features of the environment. So, for example, infants in the nonstepping window will make stepping movements when they are suspended in warm water. Stepping during the nonstepping window can also be induced by placing the infants on a treadmill. The treadmill increases leg strength by moving the leg backward and exploiting its spring-like properties. Stepping movements can also be inhibited before the start of the nonstepping window – attaching even small weights to the baby's ankles will do the trick.

These possibilities for manipulating infant stepping movements present considerable difficulties for the cortical maturation approach – since they show that stepping movements vary independently of how the cortex has developed. And they also point toward a



dynamical systems model by identifying the crucial parameters in the development of infant walking – parameters such as leg fat, muscle strength, gravity, and inertia. The brain and the rest of the central nervous system do not have a privileged position in generating this complex behavior. Instead we have a behavior that can in principle be modeled by equations tracking the interdependence of a small number of variables. Thelen and Smith have worked this idea out in great detail with a wealth of experimental studies and analyses.

Still, although walking is certainly a highly complex activity, it is not a very cognitive one. Is there support for the dynamical systems approach in a more cognitive sphere? Several examples suggest that there is. The dynamical systems approach has been profitably applied to the study of human decision-making, for example.

The Decision Field Theory developed by Jerome R. Busemeyer and James T. Townsend sets out to explain certain experimental results in behavioral economics and the psychology of reasoning in terms of the interplay of seven parameters (where agents have a choice between two actions). These parameters include settings for the strength threshold that preference need to exceed if they are to lead to action, as well as settings for the average gain from each action. A single difference equation exploits these seven parameters to fix the agent's preference at a given moment.

Another example, and one that we will look at in more detail, also derives from the work of Thelen and Smith on infant development. Thelen and Smith have developed a dynamical systems approach to how young infants understand objects.

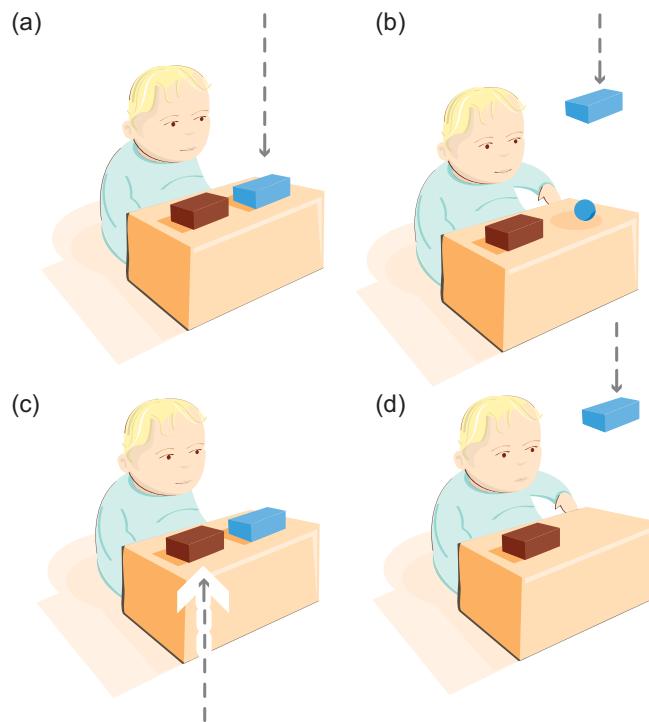
## **Dynamical Systems and the A-Not-B Error**

Object permanence is the infant's understanding that objects continue to exist when they are no longer being perceived. As we will see in much more detail in Chapter 11 (which is dedicated to object perception), object permanence emerges in stages and is intimately connected with the infant's emerging "folk physics" – with its sensitivity to the basic principles governing how physical objects behave.

One of the first to study the development of object permanence was the famous Swiss developmental psychologist Jean Piaget. In his highly influential 1954 book *The Construction of Reality in the Child* Piaget described a very interesting phenomenon.

One way to explore infants' understanding of object permanence is by looking at whether and how they search for hidden objects. Up to the age of around 7 months infants are very poor at searching for objects even immediately after they have seen them being hidden. For the very young infant, out of sight seems to be, quite literally, out of mind. From 12 months or so onward, infants search normally. But between the ages of 7 months and 12 months young infants make a striking error that Piaget termed the stage IV error and that is now generally known as the A-not-B error. Figure 6.5 illustrates a typical experiment eliciting this<sup>1</sup> error.

Infants are placed in front of two containers – A and B. They see a toy hidden in container A and reach for the toy repeatedly until they are habituated to its presence in



**Figure 6.5** The stage IV search task, which typically gives rise to the A-not-B-error in infants at around the age of 9 months. (a) The experimenter hides an object in the left-hand box. (b) The infant searches successfully. (c) But when the experimenter moves the object in full view of the infant, (d) the infant searches again at the original location. (Adapted from Bremner 1994)

container A. Then, in plain view, the experimenter hides the toy in container B. If there is a short delay between hiding and when the infants are allowed to reach, they will typically reach to container A, rather than to container B (even though they have just seen the toy hidden in container B).

Piaget himself explained the A-not-B error in terms of the infant's developing representational abilities. He suggested that it is not until they are about 12 months old that infants are able to form abstract mental representations of objects. Before then their actions are driven by sensorimotor routines. In the first stage of the task, searching for the toy in container A allows the infant to discover the spatial relationship between the toy and the container. But this knowledge only exists in the form of a sensorimotor routine. It cannot be extrapolated and applied to the new location of the toy. And so, infants simply repeat the routine behavior of reaching to container A.



### Exercise 6.3 Give in your own words Piaget's explanation of the A-not-B error.

Other cognitive and neural interpretations have been proposed. On one common interpretation, the key factor is the infant's ability to inhibit her reaching response to



**Figure 6.6** An infant sitting for an A trial (left) and standing for a B trial (right). This change in posture causes younger infants to search as 12-month-old infants do. (Courtesy L. Smith and E. Thelen)

container A. The first part of the task effectively conditions the infant to make a certain response (reaching for container A) and it is only when the infant becomes able to override that response that she can act on her knowledge of where the toy is. This ability to inhibit responses is tied to the maturation of the prefrontal cortex, which is generally held to play an important role in the executive control of behavior.

For Smith and Thelen, however, these cognitive interpretations of the A-not-B error fall foul of exactly the same sort of experimental data that posed difficulties for the cognitive interpretation of infant stepping movements. It turns out that infant performance on the task can be manipulated by changing the task. It is well known, for example, that the effect disappears if the infants are allowed to search immediately after the toy is hidden in container B. But Smith, Thelen, and other developmental psychologists produced a cluster of experiments in the 1990s identifying other parameters that had a significant effect on performance:

- Drawing infants' attention to the right side of their visual field (by tapping on a board on the far right side of the testing table, for example) significantly improves performance. Directing their attention the other way has the opposite effect.
- The most reliable predictor of infant performance is the number of times the infants reach for the toy in the preliminary A trials.
- The error can be made to disappear by changing the infant's posture – 8-month-old infants who are sitting during the initial A trials and then supported in a standing position for the B test perform at the same level as 12-month-old infants (see Figure 6.6).

If the A-not-B error were primarily a cognitive phenomenon, due either to the infants' impoverished representational repertoire or their undeveloped cortical executive system, then we would not expect infants' performance to be so variable and so easy to manipulate. It is hard to think of a cognitive/neural explanation for why standing up should make such a drastic difference, for example.

As in the infant walking case, Smith, Thelen, and their collaborators propose a dynamical systems model – the *dynamic field model*. The dynamic field represents the

space in front of the infant – the infant's visual and reaching space. High levels of activation at a specific point in the dynamic field are required for the infant to reach to that point. Thelen and Smith think about this in terms of a threshold. Movement occurs when the activation level at a particular point in the dynamic field is higher than the threshold.

Since the model is dynamical, it is critically time-sensitive. The evolution of the field has what Smith and Thelen term continual dynamics. That is, its state at any given moment depends upon its immediately preceding states. So, the activation levels evolve continuously over time. They do not jump from one state to another. What the model does is trace the evolution of activation levels in the dynamic field over time as a function of three different types of input.

- *Environmental input:* This might reflect, for example, features of the layout of the environment, such as the distance to the containers. This parameter represents the constraints the environment poses on the infant's possible actions. It will vary, for example, according to whether the infant is sitting or standing. The environmental input parameters also include the attractiveness and salience of the target, as well as contextual features of the environment, such as visual landmarks.
- *Task-specific input:* This reflects the specific demands placed upon the infant – the experimenter drawing attention to the target, for example.
- *Memory input:* The strength of this input is a function of the infant's previous reaching behavior. Since reaching behavior is partly a function of environmental input and task-specific input, the memory input reflects the history of these two types of input. And, as one would expect, it is weighted by a decay function that reflects how time diminishes memory strength.

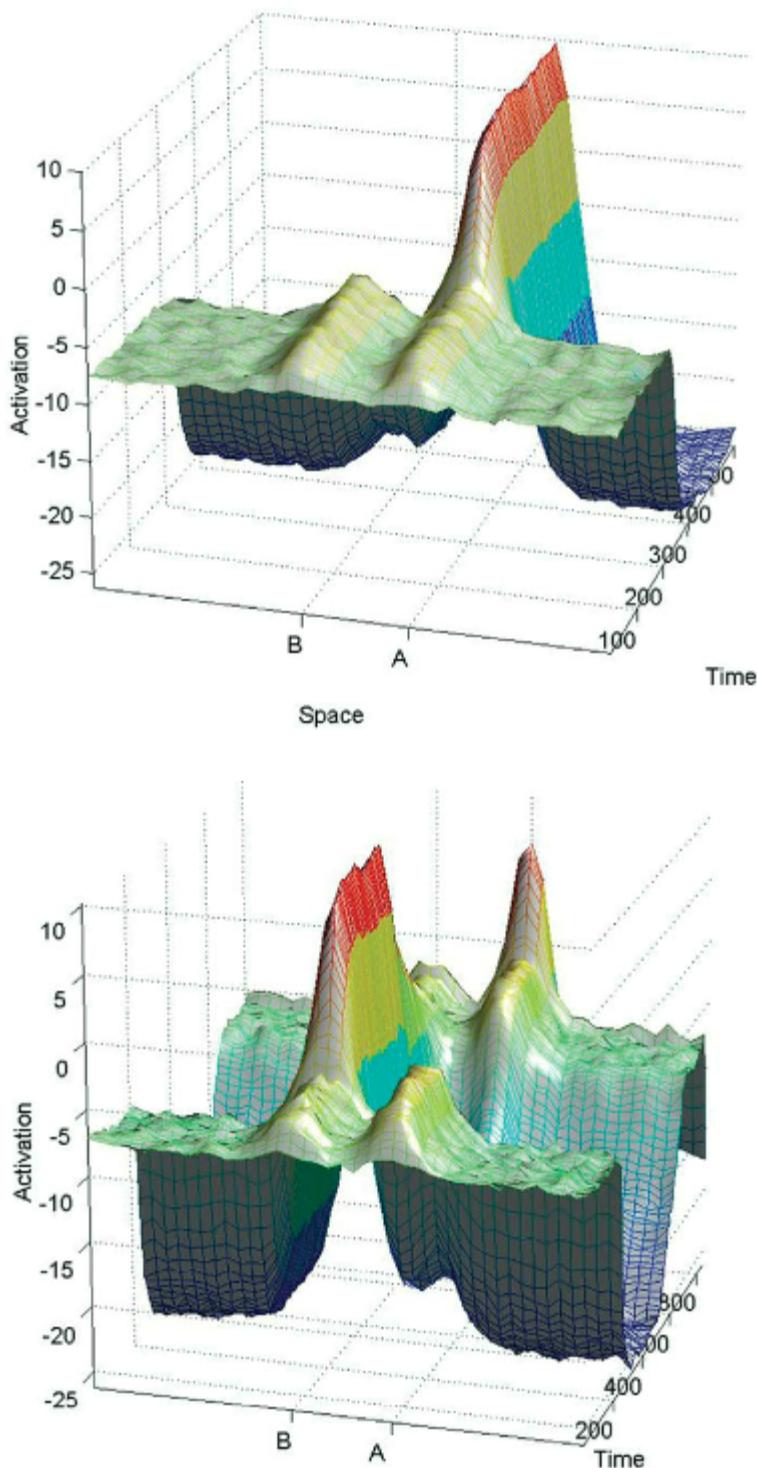
All of these parameters are coded in the same way, in terms of locations in the movement/visual field. This allows them all to contribute to raising the activation level above threshold for a specific location (either container A or container B).

And this, according to Smith and Thelen, is exactly what happens in the A-not-B error. The perseverative reaching takes place, they claim, when the strength of the memory input overwhelms the other two inputs. This is illustrated in Figure 6.7.



#### **Exercise 6.4 Explain in your own words how the dynamic field model differs from computational accounts of the A-not-B error.**

Their explanation makes no general appeal to cortical maturation, executive control, or the infant's representational capacities. And it is very sensitive to how the initial conditions are specified. If the strength of the memory input is allowed to diminish (by increasing the delay before the infant is allowed to reach, for example) then one would expect the error to diminish correspondingly – as indeed happens. The same holds for the other experimental manipulations that Smith and Thelen have uncovered. These manipulations all subtly change the inputs and parameters in the model, resulting in changes in the activation levels and hence in the infant's reaching behavior.



**Figure 6.7** Applying the dynamical field model to the A-not-B error. (a) The time evolution of activation in the planning field on the first A trial. The activation rises as the object is hidden and, owing to self-organizing properties in the field, is sustained during the delay. (b) The time evolution of activation in the planning field on the first B trial. There is heightened activation at A before the hiding event, owing to memory for prior reaches. As the object is hidden at B, activation rises at B, but as this transient event ends, owing to the memory properties of the field, activation at A declines and that at B rises.



## Assessing the Dynamical Systems Approach

The experiments and models produced by Smith, Thelen, and other dynamical systems theorists clearly give us powerful tools for studying the evolution of cognition and behavior. The explanations that they provide of the A-not-B error and how infants learn to walk seem to be both more complex and simpler than the standard type of information-processing explanations current in cognitive science. They seem more complex because they bring a wide range of factors into play that cognitive scientists had not previously taken into account, and they steer us away from explanation in terms of a single information-processing mechanism toward time-sensitive complex systems with subtle interdependencies and time-sensitivity. At the same time, their explanations seem simpler because they do not invoke representations and computations.

We started out, though, with the idea that the dynamical systems approach might be a radical alternative to some of the basic assumptions of cognitive science – and in particular to the idea that cognition essentially involves computation and information processing. Some proponents of the dynamical systems approach have certainly made some very strong claims in this direction. Van Gelder, for example, has suggested that the dynamical systems model will in time completely supplant computational models, so that traditional cognitive science will end up looking as quaint (and as fundamentally misconceived) as the computational governor.

But claims such as these ignore one of the most basic and important features of cognitive science. Cognitive science is both interdisciplinary and multilevel. The mind is too complex a phenomenon to be fully understood through a single discipline or at a single level. This applies to the dynamical systems hypothesis no less than to anything else. There is no more chance of gaining a complete picture of the mind through dynamical systems theory than there is of gaining a complete account through neurobiology, say, or AI. All of these disciplines and approaches give us deep, but partial, insights.

The contrast that Van Gelder draws between the computational governor and the Watt governor is striking and thought-provoking, but it cannot be straightforwardly transferred from engineering to cognitive science. The computational governor and the Watt governor do seem to be mutually exclusive. If we are trying to solve that particular engineering problem we need to take one approach or the other – but not both. Nothing like this holds when it comes to cognition, however. Dynamical systems models are perfectly compatible with information-processing models of cognition.

Dynamical systems models operate at a higher level of abstraction. They allow cognitive scientists to abstract away from details of information-processing mechanisms in order to study how systems evolve over time. But even when we have a model of how a cognitive system evolves over time we will still need an account of what makes it possible for the system to evolve in those ways.



Let me give an analogy. Dynamical systems theory can be applied in all sorts of areas. So, for example, traffic jams have been modeled as dynamical systems. Physicists have constructed models of traffic jams that depend upon seeing traffic jams as the result of interactions between particles in a many-particle system. These models have proved surprisingly effective at predicting phenomena such as stop-and-go traffic and the basic fact that traffic jams often occur before a road's capacity has been reached.

This certainly gives us a new way of thinking about traffic, and new predictive tools that make it easier to design roads and intersections. But no one would ever seriously propose that this new way of thinking about the collective movement of vehicles means that we no longer have to think about internal combustion engines, gasoline, spark plugs, and so on. Treating a traffic jam as an effect in a multiparticle system allows us to see patterns that we couldn't see before. This is because it gives us a set of tools for abstracting away from the physical machinery of individual vehicles. But "abstracting away from" is not the same as "replacing." Cars can be modeled as particles in a multiparticle system – but these models only make sense because we know that what are being modeled are physical objects powered (by and large) by internal combustion engines.

With this analogy in mind, look again at the dynamical field model in Figure 6.7. This model may well accurately predict the occurrence of the A-not-B error in young infants. But look at what it leaves out. It says nothing about how memory works, how the infant plans her movement, how she picks up the experimenter's cues, and soon. We don't need answers to these questions in order to construct a dynamical system model. But nor can we simply leave them unanswered. The dynamical systems approach adds a powerful tool to the cognitive scientist's tool kit, but it is unlikely ever to be the only tool.



## Summary

This chapter has explored how some cognitive scientists have used the mathematical and conceptual tools of dynamical systems theory to model cognitive skills and abilities. These models exploit the time-sensitivity that dynamical models offer in order to plot how a system evolves over time as a function of changes in a small number of system variables. We looked at two examples of dynamical systems models of child development. Dynamical systems theory offers fresh and distinctive perspectives both on how infants learn to walk and on infants' expectations about objects that they are no longer perceiving (as revealed in the so-called A-not-B error). Despite the more radical claims of some dynamical systems theorists, however, it is unlikely that traditional, information-processing models will completely disappear from cognitive science.

## Checklist

Some cognitive scientists have turned to dynamical systems theory as an alternative to traditional information-processing models of cognition.

- (1) A dynamical system is any system that evolves over time in a law-governed way, but what distinguishes the dynamical systems approach in cognitive science is the idea of studying cognitive systems with the tools of dynamical systems theory.
- (2) Dynamical models use calculus-based methods to track the evolving relationship between a small number of variables over time – a trajectory through state space.
- (3) Dynamical systems often display coupling (interdependencies between variables) and an attractor dynamics (there are points in the system's state space on which many different trajectories converge).
- (4) Cognitive systems modeled using dynamical systems theory do not display many of the classic features of information-processing systems. Dynamical models typically are not representational, computational, sequential, or homuncular.

Dynamical systems theory permits time-sensitive models of learning and skill acquisition in children.

- (1) Case studies include learning to walk in infancy, as well as performance on the A-not-B search task.
- (2) Support for the dynamical systems approach comes from experiments showing that performance can be drastically altered by manipulating factors that would typically be ignored by computational models.
- (3) The explanatory power of the dynamical systems approach does not mean that it should *replace* information-processing approaches to cognitive science.
- (4) The dynamical systems approach sheds light on cognitive systems at a particular level of organization. There is no reason to think that the level of explanation it provides should be the only one in cognitive science.

## Further Reading

Timothy Van Gelder has written a number of articles promoting the dynamical systems approach to cognitive science. See, for example, Van Gelder 1995 and 1998. The papers in Port and Van Gelder's *Mind and Motion: Explorations in the Dynamics of Cognition* (1995) contain some influential dynamically inspired studies and models (including Townsend and Busemeyer's model of decision-making), as well as theoretical statements. Thelen and Smith's 1993 edited volume *A Dynamical Systems Approach to the Development of Cognition and Action* provides more detail on their studies of infant walking, as well as contributions from other dynamical systems theorists. Their *Behavioral and Brain Sciences* article (Thelen et al. 2001) presents the model of the A-not-B error. Smith and Thelen 2003 is a more accessible introduction.

The January 2012 issue of the journal *Topics in Cognitive Science* is devoted to the complex systems approach to cognitive science, which is a branch of dynamical systems theory. For an



application of the dynamical systems approach to different areas of cognitive psychology, see Spivey 2007. For overviews and assessments of the dynamical systems approach to cognitive science, see Eliasmith 1996, Clark 1998, Clark 2001: chapter 7, Weiskopf 2004, Clearfield et al. 2009, Spencer, Thomas, and McClelland 2009, Needham and Libertus 2011, Spencer, Perone, and Buss 2011, Riley and Holden 2012, and Spencer, Austin, and Schutte 2012. For more recent reviews, see Samuelson, Jenkins, and Spencer 2015 and Spencer and Simmering 2017.





## CHAPTER SEVEN

# Bayesianism in Cognitive Science

### OVERVIEW 171

- 7.1 Bayesianism: A Primer** 172  
Degrees of Belief and Subjective Probability 173  
Conditional Probability 175  
Bayes's Rule (the Short Version) 176

- 7.2 Perception as a Bayesian Problem** 179  
The Predictive Challenge of Perception 179  
Case Study: Binocular Rivalry 182

### 7.3 Neuroeconomics: Bayes in the Brain 186

- What Is Expected Utility? 187  
Case Study: Neurons That Code for Expected Utility 190  
Probability-Detecting Neurons 193  
Utility-Detecting Neurons 194  
Combining Probability and Utility 196



## Overview

Bayesianism has become increasingly important in cognitive science. It offers a (relatively) simple model of problem solving and decision-making that has proved very profitable for modeling the mind. This chapter introduces the basic principles of Bayesianism, illustrating them through two very different case studies.

Section 7.1 lays out the elements of Bayesianism. The basic idea is that an organism's information about the world is modeled as an assignment of probabilities to different propositions about what is going on in the world. Bayesian probabilities are not like the probability that a fair coin will fall heads, or like life expectancy tables calculated from huge studies of mortality rates. Those are objective probabilities, based on frequencies. Bayesian probabilities, in contrast, are subjective. They reflect an organism's best guess. That best guess is updated as new information comes in. Bayesians propose that this updating takes place according to a (relatively) simple rule called Bayes's Rule, named after Thomas Bayes, an eighteenth-century clergyman and statistician.

Section 7.2 shows how perception can be seen as a Bayesian problem. Perceptual systems have to work backward from noisy data to a consistent, coherent, and, it is hoped, accurate, model of how the world is. This is an inference problem. The brain has to make a reverse inference from the



noisy data to a hypothesis about the layout of the world. This is exactly the kind of inference for which Bayes's Rule is particularly well suited. We will illustrate this through a Bayesian model of the phenomenon of *binocular rivalry*.

The theory of expected utility is an extension of Bayesian principles to decision-making. Section 7.3 introduces the principle of expected utility and explains how expected utility is calculated. Our second illustration of Bayesian modeling comes from the growing field of *neuroeconomics*, which uses economic tools such as the theory of expected utility to studying the brain. We will look at a series of experiments tracking individual neurons in an area of the parietal cortex known as the lateral intraparietal area (usually abbreviated as LIP). There appear to be neurons in LIP that code for probability and for analogs of utility and expected utility.

## 7.1

### Bayesianism: A Primer

Remarkably, given how powerful it is as a modeling tool, the basic elements of Bayesianism are really pretty straightforward. There are three key ideas.

- Belief comes in degrees.
- Degrees of belief can be modeled as probabilities, and so have to conform to the basic principles of the probability calculus.
- Learning takes place by updating probabilities according to Bayes's Rule.

This section presents these three key ideas in turn. From a technical point of view, a basic knowledge of the probability calculus is really all you need to see what is going on.



**Figure 7.1** An illustration purporting to be of Thomas Bayes from a 1936 book on the history of life insurance. (From Wikimedia Commons)



## Degrees of Belief and Subjective Probability

It is natural to talk about cognitive systems having beliefs about their environment. But what are beliefs? Most people think about beliefs in a very particular way. They think of them as two-valued states. A belief is either true or it is false. If you believe that there is a predator hiding in the woods, then either there is a predator there and your belief is correct, or there is no predator and your belief is incorrect. There is no in-between state. Beliefs are, as it were, pass/fail.

For Bayesians, this is completely the wrong approach. They think that belief is not an on-off state. It comes in degrees. At one end of the spectrum are things in which you are completely confident. At the other end of the spectrum are things that you have absolutely no confidence in at all. From a Bayesian perspective, the interesting things all happen in between these two extremes. Most of what we believe about the world we have some confidence in, but not complete confidence.



**Exercise 7.1** Give examples of (a) something in which you have complete confidence; (b) something in which you have no confidence whatsoever; (c) something in which you are fairly confident, but not completely so; and (d) something in which you have a slight degree of confidence.

The second key idea of Bayesianism is that degrees of belief are probability assignments. Bayesians replace the vague idea of having more or less confidence in some proposition with the much more precise notion of assigning a particular numerical probability to that proposition. So, rather than describe an organism as being more confident than not, but still some way short of completely confident, that there is a predator in the woods, a Bayesian might say that the organism assigns a probability of 0.7 to there being a predator in the woods.



**Exercise 7.2** Go back to Exercise 7.1 and assign numerical probabilities to your answers to (a) through (d).

It is essential to Bayesians that degrees of belief obey the fundamental principles of the probability calculus. So, we can use the rules of the probability calculus to update and combine our degrees of belief. Box 7.1 offers a quick refresher of the basic rules of probability, but here are some examples that will probably seem familiar.

- If we assign probability  $p$  to some sentence  $S$ , then we have to assign probability  $1 - p$  to **not-S**, the negation of  $S$ .
- If we assign probability  $p$  to  $S$  and probability  $q$  to  $R$ , and we know that  $S$  and  $R$  are (probabilistically) independent of each other, then we have to assign probability  $p \times q$  to **S AND R** (i.e., to  $S$  and  $R$  both holding).
- If we assign probability  $p$  to  $S$  and probability  $q$  to  $R$ , and we know that  $S$  and  $R$  are mutually exclusive, then we have to assign probability  $p + q$  to **S OR R** (i.e., to at least one of  $S$  and  $R$  holding).

## BOX 7.1 Basic of the Probability Calculus

There are different but equivalent ways of presenting the probability calculus. Probabilities are sometimes assigned to sentences and sometimes to events. If you think in terms of sentences, then it is most natural to formulate the probability calculus in logical terms (talking about conjunctions and disjunctions, for example). If you think in terms of events, then it is most natural to use elementary set theory (talking about intersections and unions of events, for example, I'll use the framework of sentences).

There are four basic principles defining how probabilities behave:

- Basic principle 1: Probabilities are numbers between 0 and 1
- Basic principle 2: All impossible sentences have probability 0
- Basic principle 3: All necessary truths (such as " $2 + 2 = 4$ ") have probability 1
- Basic principle 4: If sentences  $P$  and  $Q$  are logically equivalent, then  $p(P) = p(Q)$

Then, with these basic principles in place, the probability calculus contains simple rules that tell us how to assign probability to complex sentences built up from simpler sentences for which we know the probabilities.

### The Negation Rule

If sentence  $S$  has probability  $p$ , then its negation **not-S** ( $\neg S$ ) has probability  $1 - p$

### The Disjunction Rule (Restricted)

If sentences  $R$  and  $S$  are mutually exclusive, then the probability of **R or S** is  $p(R) + p(S)$

### The Conjunction Rule (Restricted)

If sentences  $R$  and  $S$  are independent of each other (i.e., the presence of one does not make the other more likely), then the probability of **R and S** is  $p(R) \times p(S)$

In order to lift the restrictions, we need to apply the concept of **conditional probability**, as explained in the text. The conditional probability of  $S$  conditional upon  $R$  (written  $p(S|R)$ ) is the probability that  $S$  holds, on the assumption that  $R$  holds. So, for example, the probability of throwing a 4 with a 6-sided die is  $1/6$ . But the probability of throwing a 4, conditional upon throwing an even number is  $1/3$ .

We can use conditional probability to define

### The Conjunction Rule (General)

$$p(R \text{ and } S) = p(R|S) \times p(S)$$

And then, with this general definition of conjunction, we can define

### The Disjunction Rule (General)

$$P(R \text{ or } S) = p(R) + p(S) - p(R \text{ and } S).$$

That's it!



From a Bayesian perspective, though, the most important inference rule for reasoning with degrees of belief is Bayes's Rule. But to understand Bayes's Rule, we first need to understand the concept of conditional probability.

## Conditional Probability

Here is the basic idea of conditional probability. The probability of A, conditional upon B, is the probability that A holds, relative to the assumption that B holds. In other words, you assume for the moment that B holds and then calculate the probability of A on that assumption. We write the probability of A, conditional upon B, like this:

$$p(A|B)$$

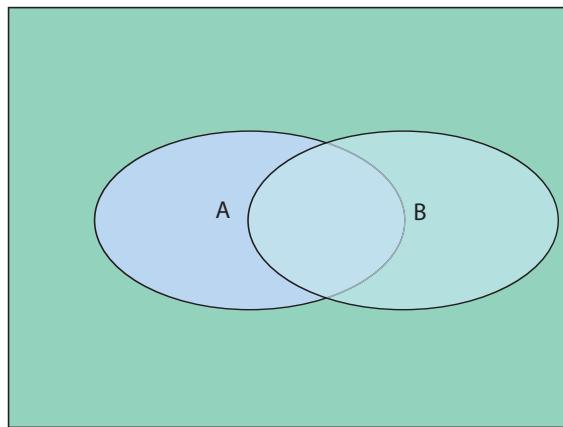
So, if "A" stands for "There is thunder" and "B" for "It is raining," then  $p(A|B)$  is the probability that there is thunder, if we assume that it is raining.



**Exercise 7.3** Explain why, if "A" stands for "There is thunder" and "B" for "It is raining," then you would expect  $p(A|B)$  to be greater than  $p(A)$ .

Here is a diagram that will get us started on understanding conditional probability.

You can think about probabilities in terms of a sample space. The sample space is, as it were, all the possibilities that there are. Individual probabilities are regions of the sample space, which we measure by how much of the sample space they occupy. So, for example, in the diagram the probability of A is the proportion of the sample space in which A is true. If A is true in 30 percent of the sample space, then the probability of A is 0.3. And the same holds for B, of course.



**Figure 7.2** A diagram showing (a) the proportion of the probability space in which A is true, (b) the proportion of the probability space in which B is true, and (c) the intersection of A and B (which is the region where A and B are both true).

Now, what about the probability of the conjunction A&B? Well, it's easy to see what that is from the diagram. The probability  $p(A \& B)$  is the proportion of the sample space in which both A and B are true – where the A-space overlaps with the B-space.

So, finally, what about the probability of A conditional upon B? Well, for  $p(A/B)$  we also need to consider the region of the sample space that is both A and B. This is the intersection of A and B. But we are not looking at how that region of the sample space relates to the sample space as a whole. Since we are looking at A, conditional upon B, we are interested only in the relation between the A and B region of the space and the B region of the space. In other words,  $p(A/B)$  is the proportion of the intersection of A and B to the B-space.

Well, if you find that convincing, then you will have no problem with the formal definition of conditional probability, because the formula for conditional probability is basically a translation of what I've just said into the language of probability. Here it is:

$$p(A/B) = \frac{p(A \& B)}{p(B)}$$

In other words, to derive the conditional probability  $p(A/B)$  you take the probability of A and B both holding (i.e.,  $p(A \& B)$ ) and divide by the probability that B holds.



**Exercise 7.4** Suppose you have a fair coin, which you toss twice. Let "A" stand for "the first toss comes up heads" and "B" for "the second toss comes up heads." Calculate (a)  $p(A)$ ; (b)  $p(B)$ ; and (c)  $p(A \& B)$ . Then (d) use the formula in the text to calculate  $p(A/B)$  – i.e., the probability that the second toss comes up heads, conditional on the first toss coming up heads.

## Bayes's Rule (the Short Version)

This section offers a nontechnical presentation of the basic idea behind Bayes's Rule. The rule actually follows fairly straightforwardly from the definition of conditional probability. See Box 7.2 for the long version, with more details.

To see what Bayes's Rule allows us to do, it helps to think about situations in which you have a hypothesis that you are trying to decide whether to accept or reject. You have some evidence for the hypothesis. So, the obvious question to ask is: How strong is that evidence?

When we ask that question, what we are really asking is: How likely is it that the hypothesis is true, given the evidence? And that, in turn, is really a question about conditional probability. We are asking about the probability of the hypothesis conditional upon the evidence:  $p(\text{Hypothesis}/\text{evidence})$ .

Let's make things a little simpler with some obvious abbreviations. I'll use "H" for the hypothesis and "E" for the evidence. So, what we are trying to discover is  $p(H/E)$ .

Often, in this sort of situation, we have information about how likely the evidence is, given the hypothesis. In other words, we might know the conditional probability  $p(E/H)$ .



## BOX 7.2 Deriving Bayes's Rule

Here's the formula for Bayes's Rule, using the abbreviations in the main text ("E" for evidence and "H" for hypothesis):

$$p(H/E) = \frac{p(E/H)p(H)}{p(E/H)p(H) + p(E/\neg H)p(\neg H)}$$

Let's start with the top of the equation (the numerator). This is the probability of the evidence, conditional upon the hypothesis, multiplied by the probability of the hypothesis. The key thing to understand here is that this amounts to the probability of the evidence and the hypothesis both holding. In other words,  $p(E/H)p(H)$  gives us  $p(E \ \& \ H)$ .

Now, look at the bottom of the equation (the denominator). This has two parts, connected by an addition sign. The first part is the same as the numerator. It is essentially  $p(E \ \& \ H)$ . So, it is the probability of the evidence and the hypothesis both holding. The second part works the same way, so that  $p(E/\neg H)p(\neg H)$  comes out as  $p(E \ \& \ \neg H)$ . This is the probability that the evidence holds, but not the hypothesis. In other words, that you have a false positive.

When you put everything together in the denominator, you should see that it gives the probability that **either** you have a true positive test (the hypothesis and the evidence both hold) **or** you have a false positive test (where the evidence holds but not the hypothesis).

But these are the only two options. Since you've had a positive test, it has to be either true or false. So, what the denominator really adds up to, then, is simply the probability of getting a positive test in the first place. In other words, you can write the denominator much more simply as  $p(E)$ .

So really, then, Bayes's Rule is a lot simpler than it initially seems. Here's how we can simplify it:

$$(1) \quad p(H/E) = \frac{p(E/H)p(H)}{p(E/H)p(H) + p(E/\neg H)p(\neg H)}$$

which simplifies to

$$(2) \quad p(H/E) = \frac{p(E\&H)}{p(E\&H) + p(E\&\neg H)}$$

which simplifies to

$$(3) \quad p(H/E) = \frac{p(E\&H)}{p(E)}$$

But now – look where we have arrived. This is a straightforward application of the definition of conditional probability.



**Exercise 7.5** To make sure that you understand the reasoning here, go through the reverse of the simplification process, starting with the definition of conditional probability and ending up with Bayes's Rule.

In that situation, what we need to do is to find a way of “inverting” this probability, so that we can get from  $p(E/H)$  to  $p(H/E)$ .

Imagine a medical scenario. You test positive for a nasty disease. So, the hypothesis ( $H$ ) is that you actually have the disease. The evidence ( $E$ ) is the positive test. And suppose you know how reliable the test is. Then you know how likely it is that you will test positive if you actually do have the disease. In other words, you know  $p(E/H)$ . But what you really want to know is  $p(H/E)$  – that is, how likely it is that you have the disease, given that you have tested positive.

Here is some terminology that tends to appear in discussions of Bayesian approaches to cognitive science.

### **Posterior probability**

This is the probability that you end up with, after applying Bayes’s Rule. It is  $p(H/E)$ .

### **Prior probability**

This is the probability that you originally assign to the hypothesis. It is  $p(H)$ .

### **Likelihood of the evidence**

This is the probability that you’ll get the evidence, if the hypothesis is true. It is  $p(E/H)$ .

To continue with the medical example, the likelihood of the evidence is the reliability of the test. If it accurately detects 99 percent of cases of the disease, then the likelihood of the evidence is 0.99. The prior probability would be given by the frequency of the disease in the population. So, if the disease afflicts 1 in 10,000 people then the prior probability would be 0.0001

So, using this terminology we can write Bayes’s Rule in words like this

$$\text{Posterior probability of the hypothesis} = \frac{\text{Likelihood of the evidence} \times \text{Prior probability of the hypothesis}}{\text{Probability of the evidence}}$$

If you remember this formulation in words, then you’ll understand the basic conceptual foundation of Bayesian updating. If you can also remember the formula in the probability calculus, as given in Box 7.2, then you will be able to plug numbers in and apply Bayes’s Rule to solve specific problems.

Either way, however, the basic idea of Bayesianism is that Bayesian agents update their beliefs by applying Bayes’s Rule. If you are a Bayesian agent, then you start off with a set of prior probabilities. These are the probabilities that you assign to the different hypotheses about how things might turn out. You are also aware of how likely it is, for each of these hypotheses, that you will encounter different forms of evidence. So, you know various likelihood probabilities. And then, as the evidence comes in, you apply Bayes’s Rule to derive your posterior probabilities.



It's important to remember all of the different things that go into Bayes's Rule. It is easy, for example, to underestimate the importance of the priors. But that can be a grave mistake, as you can see by looking at the equation again. The lower the prior probability of the hypothesis is, then the lower the numerator will be in the equation (i.e., the part above the line). Since the denominator (the part below the line) is less than 1, that means that the posterior probability will end up being lower.

This means that, if the prior probability is very low, then the posterior will also be very low, no matter what the other values are. This is very important in the context of medical diagnosis. Suppose again that you test positive for a nasty disease that is extremely rare. Even if the test is highly reliable (say, 0.9999), then it is highly unlikely that you actually have the disease.

Here's why. Suppose that only 1 person in every 100,000 has the disease. Now imagine that 100,000 are tested. The test is highly reliable, so we can assume that that person will test positive. But, even though the test is highly reliable, it will still misdiagnose one person in every 10,000. So, there will be around ten false positives in our population of 100,000 people. That means that there will be eleven people testing positive, only one of whom actually has the disease. So, your odds are actually not that bad!



**Exercise 7.6** Work through the example in the text, but instead of the prior probability being 1 in 10,000, take it to be 1 in 100. How likely is it now that you have the disease?

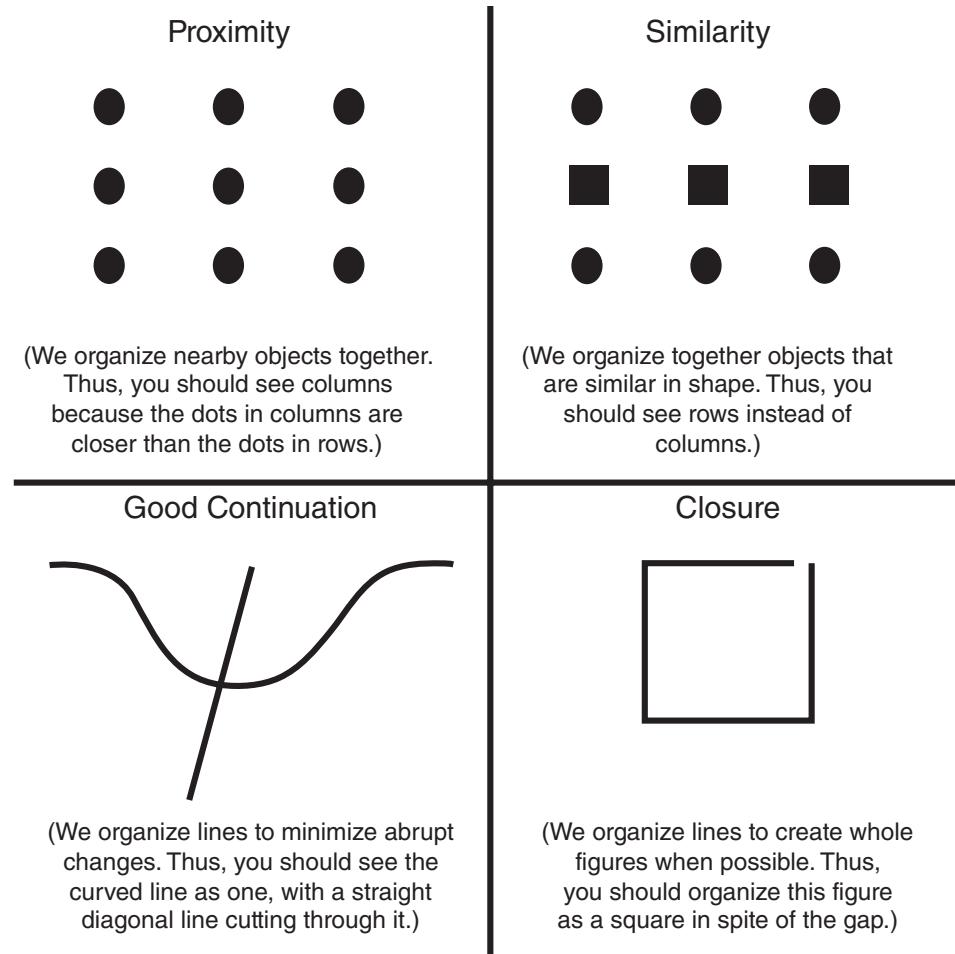
As this little example shows, Bayes's Rule is very powerful, despite its simplicity, and it can give surprising results. In the next section we'll look at how it can be applied in the context of perception.

## 7.2 Perception as a Bayesian Problem

After working through the theory and the math, it is time to see how it can be applied. This section develops a case study of how Bayes's Rule might be used in cognitive science modeling. We will be looking at visual perception – and in particular at a puzzling phenomenon known as *binocular rivalry*. We'll start by looking at why perception seems a good candidate for Bayesian approaches. And then we'll turn to binocular rivalry, and see how we can make sense of it by assuming that the visual system is engaged in Bayesian updating using Bayes's Rule.

### The Predictive Challenge of Perception

Perception is an obvious place to apply Bayesian ideas. Our perceptual systems deliver a model of the environment that is, by and large, fairly accurate – at least for practical purposes. But this model is grossly underdetermined by the information that actually reaches the sensory systems. So, how do our perceptual systems get from proximal sensory stimulation to a full-blown model of the distal environment?



**Figure 7.3** Four of the seven Gestalt principles of grouping, illustrated and explained.  
(Downloaded from [www.skidmore.edu/~hfoley/PercLabs/Shape.htm](http://www.skidmore.edu/~hfoley/PercLabs/Shape.htm))

Back in the second half of the nineteenth century, the German physicist and psychological pioneer Hermann von Helmholtz suggested, very plausibly, that perception was at bottom a process of unconscious inference. Inferential models of perception were also an important part of the New Look perceptual psychology pioneered by Jerome Bruner in the 1940s. Some decades later, in 1983, Irving Rock published a famous book entitled *The Logic of Perception*, which offered a very distinctive version of the inferential model of perception.

Rock was very much influenced by the Gestalt school of perceptual psychology, including figures such as Max Wertheimer, Kurt Koffka, and Wolfgang Köhler (in Germany), and Rudolf Arnheim (in the USA). He took from the Gestalt psychologists the idea that perceptual inference is largely top-down and holistic. From the Gestalt perspective, what the visual system does is impose structure on an essentially unstructured retinal image. And in doing this it uses general principles about how objects and groups of objects form organized patterns. These are the famous Gestalt principles of grouping. Four of these principles are illustrated in Figure 7.3.



So, from a Gestalt perspective and on inferentialist views such as Rock's, the perceptual systems form hypotheses about the distal environment based on general principles of organization and grouping.

But still, you might reasonably wonder how exactly those Gestalt principles work. How do those general principles of organization and grouping actually operate to structure how we see the world?

Bayesian models of perception have got an answer to this question. They think that perceptual systems make probabilistic inferences, so that perceptual systems exploit a constantly updated body of probabilistic knowledge. This probabilistic knowledge can be described in terms of the basic Bayesian concepts introduced in the last section.

- Perceptual systems are aiming to derive a hypothesis ( $H$ ) about the layout of the distal environment.
- Ultimately all they have to go on is the evidence ( $E$ ) provided by sensory stimulation at the retina, or at the membrane window of the cochlea.
- Each perceptual system aims for the hypothesis that is most probable given the evidence. So, what it is ultimately interested in are conditional probabilities of the  $p(H/E)$  variety. These are the posterior probabilities.
- Perceptual systems store information about the likelihood of different environmental set-ups. These are the prior probabilities –  $p(H)$ .
- Perceptual systems also store information about how likely different types of sensory stimulation are, given different layouts of the distal environment. These are likelihoods, conditional probabilities of the form  $p(E/H)$ .

Perhaps you can see where the Gestalt principles might fit into this overall picture?

The Gestalt principles are essentially principles about the probable structure of the environment. They are principles that govern how it might be reasonable to work backward from patterns in the retinal image to the objects from which those patterns ultimately originate. So, it is helpful to think of them as principles about the probability of different types of environment setup. Or in other words, they function as Bayesian priors. Look again at Figure 7.3 to appreciate this. The bottom left-hand box illustrates the Principle of Continuation. Effectively, what this says is that a hypothesis about the layout of objects in the environment that contains abrupt changes should have a lower prior probability than one with fewer or no abrupt changes.

Looked at in this way, Gestalt principles are examples of Bayesian priors. They are fundamental elements in the process of probabilistic inference that perceptual systems use to formulate hypotheses about how the distal environment is laid out. So, the next question is: How do perceptual systems make probabilistic inferences that end up in a model of how things are in the external environment?

There are no prizes for guessing that Bayesians think that this is all done using Bayes's Rule! Let's look now at an example of how this might work.



## Case Study: Binocular Rivalry

Our first case study of how Bayesian approaches can be applied in cognitive science looks at an intriguing phenomenon from visual perception known as *binocular rivalry*.

The Italian scientist and playwright Giambattista della Porta discovered in the late sixteenth century that when separate images are presented to each eye, what you actually perceive alternates between the two images. The image that you see switches seemingly at random. Figure 7.4 illustrates two examples of stimuli that can be used to elicit the phenomenon. You can check the effect out for yourself. There are experimental demonstrations on the internet. See the Further Reading section and website for details and directions.

Binocular rivalry is really a special case of the more general phenomenon of perceptual rivalry. If you look at an ambiguous figure, such as one of the examples in Figure 7.5, then your visual system will lock onto one of the available interpretations and you have to work hard to start seeing it the other way – to switch from seeing a duck to seeing a rabbit for example.

Perceptual rivalry in general, and binocular rivalry in particular, are very interesting for a number of reasons. The key point is that what you perceive changes, even though the stimulus remains the same. In binocular rivalry, there does not seem to be anything about the image in front of either eye that prompts the switch. So, whatever the explanation is for the alternating perceptions, it must lie either within the visual system, or downstream in more central processing.

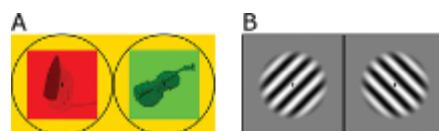
Different explanations have been proposed for binocular rivalry. Some early theorists, including della Porta himself, took it as evidence for the view that we only ever see with one eye at a time. Others speculated that the switch might be due to attentional factors. Much more recently, though, a simple and elegant Bayesian explanation has been proposed by Jakob Hohwy, Andreas Roepstorff, and Karl Friston in an article published in the journal *Cognition* in 2008.

To appreciate their basic point, think about the challenge that the visual system faces in a typical binocular rivalry situation, like that depicted in the left-hand pair of stimuli in Figure 7.4. The visual system has to decide what it is looking at. Does the distal environment contain a red iron? Or a green violin? Or some sort of mish-mash composite object that is part iron, part violin, colored both red and green? (To simplify, I will gloss over the fact that we are really dealing with a picture of a red iron, rather than a red iron.)

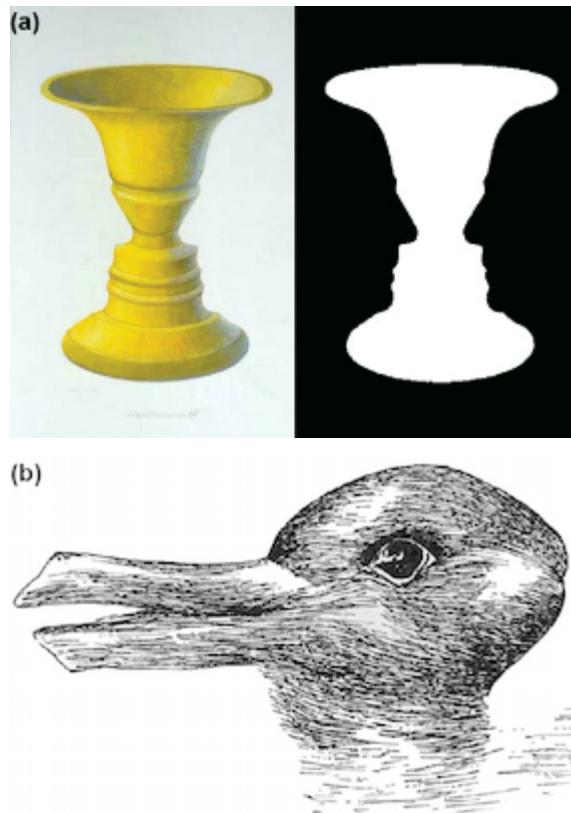
In Bayesian terms, there are three different hypotheses:

Hypothesis 1 (H1): Red iron

Hypothesis 2 (H2): Green violin



**Figure 7.4** Two examples of stimuli used to elicit binocular rivalry. (Figure 1 from Freyberg et al. 2015)



**Figure 7.5** Two well-known ambiguous figures: (a) Rubin's vase and (b) the duck–rabbit illusion. (Downloaded from <https://upload.wikimedia.org/wikipedia/commons/b/b5/Rubin2.jpg> and [https://commons.wikimedia.org/wiki/File:Duck-Rabbit\\_illusion.jpg](https://commons.wikimedia.org/wiki/File:Duck-Rabbit_illusion.jpg), respectively)

Hypothesis 3 (H3): Red-green iron-violin composite object

So – the visual system's job is to decide which hypothesis to accept, given the evidence it has. And the evidence is the images that are presented to each eye. In other words:

Evidence (E): Picture of red iron (L eye) and picture of green violin (R eye)

Translating this all back into the vocabulary that we looked at in Section 7.1, what the visual system needs to figure out is which of the following conditional probabilities is the highest:

$$p(H1/E)$$

$$p(H2/E)$$

$$p(H3/E)$$

These are the *posterior probabilities* – the probability, for each hypothesis, that it is true, given the available evidence. A Bayesian visual system will accept the hypothesis with the highest posterior probability.

This is where Bayes's Rule comes into play. The visual system can apply Bayes's Rule to derive the posterior probability for each hypothesis, and then it just needs to compare the resulting posterior probabilities.

In order to apply Bayes's Rule, two more types of information are needed. The first is the *likelihoods*. Recall that these are conditional probabilities, specifying how likely the relevant evidence would be if the corresponding hypothesis were true. So, for example, how likely it would be to have a red iron image on the left and a green violin image on the right if what was out there was a red iron. There are three likelihoods to consider:

$$p(E/H1)$$

$$p(E/H2)$$

$$p(E/H3)$$

The visual system also needs to consider the *prior probability* of each hypothesis. How likely is it that what's in front of it is a red iron (H1), a green violin (H2), or a red-green iron-violin composite object (H3)? So, that gives us three more probabilities to factor in:

$$p(H1)$$

$$p(H2)$$

$$p(H3)$$

Bayes's Rule tells us how to use the priors and the likelihoods to calculate the posterior probabilities.

But wait, you might ask: How on earth is the visual system going to be able to assign numerical values to all of these conditional and unconditional probabilities? It seems implausible to think that the visual system might assign a specific numerical value, say 0.47, to the prior probability of the hypothesis (H1) that what's out there is a red iron – or to the likelihood  $p(E/H3)$  of there being a red iron image on the left and a green violin image on the right if what was out there was a composite red-green iron-violin?

This is a very good question. It is certainly part of strict Bayesianism that a rational Bayesian agent will always be able to assign numerically definite probabilities to any possible outcome it considers. And opponents of Bayesianism often object that this is not a feasible or desirable requirement. From a cognitive science perspective, though, it is important to realize that you can apply Bayesian techniques and models to the visual system without assuming that the visual system assigns specific numbers to the priors and likelihoods.

To see how this works, let's look again at the formula for Bayes's Rule. Here it is:

$$\text{Posterior probability of the hypothesis} = \frac{\text{Likelihood of the evidence} \times \text{Prior probability of the hypothesis}}{\text{Probability of the evidence}}$$



The visual system is trying to compare the three hypotheses: H1, H2, and H3. The denominator (the bottom of the equation) is going to be the same for each hypothesis. (If you want to check this, look in Box 7.2 at how  $p(E)$  is calculated.)

So, for the purposes of comparing H1 through H3, the visual system can effectively ignore the denominator. Really, then, what it needs to do is calculate, for each hypothesis, the product of the likelihood of the evidence and the prior probability. For H1, for example, it would need to multiply the likelihood  $p(E/H1)$  by the prior  $p(H1)$ . And so on for H2 and H3.

But now suppose we make a simplifying assumption. Suppose we think that all the likelihoods are the same. That means that the visual system thinks that the evidence (an image of a red iron on the left and a green violin on the right) would be equally likely whether there was a red iron, a green violin, or a red-green iron-violin in front of it. Then the visual system can ignore the likelihoods too, and all it needs to look at are the prior probabilities.

Looking at the priors, it seems reasonable to think that there is no particular reason for the visual system to think that it is any more or any less probable that there would be a red iron in front of it than a green violin. So, you might think that in this case

$$p(H1) = p(H2).$$

But at the same time, the visual system will probably think it very unlikely that it will run into a red-green iron-violin. So, it would assign a much lower prior probability to H3. Hence:

$$p(H3) < p(H1) \text{ and } p(H3) < p(H2).$$

So, we can put this all together to get the key idea in Hohwy, Roepstorff, and Friston's 2008 paper. The end result is that there are joint winners. The visual system assigns the same posterior probabilities to H1 and to H2, with both clearly preferred to H3. That gives:

$$p(H1/E) = p(H2/E).$$

But now the visual system has nothing to go on to choose between H1 and H2. It could be a red iron, or it could be a green violin. So, what does it do?

Well, unable to decide between H1 and H2, the visual system switches between them, more or less at random. It doesn't try to construct some sort of composite perceptual image, because that would correspond to a much less likely outcome. Put another way, the prior probability of the environment containing a composite object is much lower than the prior probability that it contains an iron or the prior probability that it contains a violin.

According to Hohwy, Roepstorff, and Friston, this creates the characteristic effect of binocular rivalry, as a rational response to a process of Bayesian updating.

Notice that all this was done without assigning any precise numerical probabilities. Admittedly, I made it all much simpler with the assumption that the likelihoods were the same for all three hypotheses. But actually, lifting this assumption doesn't make it much more complicated.

It certainly seems reasonable that the likelihoods of the evidence would be the same for H1 and H2. In each case, the hypothesis only explains part of what is going on in the evidence. If H1 holds, then that explains why there is a red iron image on the left. But the green violin is a complete mystery. But the situation is completely symmetrical for H2. If H2 holds then that explains why there is a green violin image on the right. But now the red iron image is completely mysterious. Either way you have to ignore part of the stimulus so that the rest of it makes sense.

So, consider the likelihood  $p(E/H3)$ . This is the likelihood that you would get the evidence (a red iron image on the left and a green violin image on the right) if what were out in the world were a composite red-green violin-iron. How likely is that? Well, probably not very likely at all. If the visual system is even remotely reliable, then you would expect one or both eyes to generate an image of what is out in the world, namely, a red-green iron-violin. It seems unlikely that a composite object would yield two distinct retinal images, neither of which really corresponds to what is supposed to generate them.

So, the visual system will assign a lower probability to  $p(E/H3)$  than to either  $p(E/H1)$  or  $p(E/H2)$ . And then, just as before, Bayes's Rule will continue to have the posterior probabilities of H1 and H2 come out the same, conditional upon the evidence. And each of them will be more probable than H3, given the available evidence. The pieces are still in place for a Bayesian visual system to display the binocular rivalry effect.



**Exercise 7.7** Write down in your own words a summary of the Bayesian explanation of binocular rivalry.

## 7.3 Neuroeconomics: Bayes in the Brain

We turn now to a second, and very different, example of how Bayesian principles can be fruitfully used to study the brain. This example uses another important dimension of Bayesianism – the theory of expected utility. It comes from neuroeconomics.

Neuroeconomics is an interdisciplinary area within cognitive science. It is located at the interface between neuroscience, on the one hand, and microeconomics, the psychology of reasoning, behavioral finance, and decision theory, on the other.

- *Microeconomics* is the branch of economics that studies how individual consumers allocate resources (in contrast to *macroeconomics*, which studies the economy as a whole).
- *The psychology of reasoning* is the experimental study of how people reason and make decisions.
- *Behavioral finance* focuses on how individuals make investments, and what that can tell us about the financial markets.
- *Decision theory* is the mathematical theory of rational choice and decision-making. Bayesian approaches are very influential in contemporary decision theory.



In brief, neuroeconomics studies how brains deal with money, investments, and risky choices. It uses the tools of neuroscience, obviously, but combines them with experimental paradigms that psychologists have developed to study reasoning and data from financial markets, as well as theoretical models from decision theory and elsewhere.

Neuroeconomics has, broadly speaking, two different dimensions. One dimension is studying the neuroscience of decision-making – looking to see how value is computed in the brain and how such computations play into decision-making. The second dimension works in the opposite direction, as it were. It takes the theoretical tools that Bayesianism uses to study decision-making and then applies those tools to study the brain. This is the dimension of neuroeconomics that we will be exploring.

The dominant theoretical model in neuroeconomics is Bayesian expected utility theory. This is a development of the basic Bayesian approach developed and discussed in earlier sections. We looked at Bayesian approaches to theoretical reasoning, focusing on how to measure the support evidence provides for a hypothesis. The principal tool was the theory of probability, particularly Bayes's Rule.

Bayesian expected utility theory also has the theory of probability at its core, but it goes beyond theoretical to incorporate practical reasoning. To appreciate the distinction, consider the difference between these two questions:

What should I believe? (theoretical reasoning)

What should I do? (practical reasoning)

Probability theory and Bayes's Rule are enough to answer the theoretical question (according to strict Bayesians). But we need more machinery for the practical question.

That's where expected utility theory comes into play. The concept of utility allows us to model how much an individual values different outcomes. Being able to measure value in this way allows us to choose between different courses of action as a function of two things. We need, first, to take into account, for each available course of action, how much we value each of its different possible outcomes. And then, second, we need to consider how likely each of those outcomes is.

Combining probability and utility allows a rational chooser to maximize expected utility. We'll look in more detail at how this works next. And then we'll look at some exciting experiments using single-cell recordings on monkeys that have identified individual neurons in the parietal cortex that seem to code for analogs of the basic concepts of expected utility theory.

## What Is Expected Utility?

The theory of expected utility is the cornerstone of many social sciences, but it is particularly prominent in economics. Although the actual concept of utility did not appear until the nineteenth century, it has its roots in the sixteenth and seventeenth centuries when mathematicians (and professional gamblers) were starting to develop ideas about probability in order to understand games of chance.

The best place to start is with the concept of expected value (and, by the way, what we are talking about here is monetary value, not artistic value, or esthetic value). Here's a simple example. What is the most that you would pay for the opportunity to play a game in which you win \$10 if a fair coin is tossed and lands heads, but get nothing if it comes up tails?

Let's assume, for the sake of argument, that you are what is called risk-neutral. That means that you have no views either way about taking risks. You are not prepared to pay extra because you enjoy the thrill of gambling. And nor do you require any special compensation to bet money on the toss of a coin.

It seems intuitive that a risk-neutral person would pay up to \$5 for the chance to play this game. Why? Because there are two possible outcomes – heads or tails. If the coin comes up heads, she will walk away with \$10. If it comes up tails, she will get \$0. So, her average return will be \$5. If she pays less than that, she might reasonably think she's got a bargain. But if she pays more, then it looks as if she has been fleeced (unless she is risk-loving and prepared to pay a premium because she loves gambling). So, in the standard terminology, the *expected monetary value* of this game is \$5.

You can see why, from the perspective of a seventeenth-century professional gambler, expected value would be a very useful concept. If you know the expected value of a gamble and your opponents do not, then you can set the odds so that you are guaranteed to come out on top over the long term. In essence, this is why casinos and lotteries always win in the long run. They set the odds so that punters are always paying more than the expected value of the game or the ticket.

But in some instances the concept of expected value gives strange and implausible results. Suppose, for example, that you have the opportunity to play the following game.

#### *St. Petersburg Game*

A fair coin is tossed. If the coin lands heads, you receive \$2 and the coin is tossed again. If it lands tails, the game ends. But if it lands heads, then you will receive \$4 and the coin is tossed again. The game will continue as long as the coin lands heads. At each round the pay-off (for heads) will be twice what it was in the previous round.

Think about the expected value of the St. Petersburg game. The expected value of the game is the expected value of continuing to get heads. On each toss of a fair coin, the probability of heads is 0.5. So the expected value of the game is  $(0.5 \times \$2) + (0.5 \times \$4) + (0.5 \times \$8) + \dots$  and so on indefinitely.

You should be able to see that the expected value is infinite (if not, try rewriting the sum as  $1 + 2 + 4 + 8 + \dots$ ). So, if we think about this the same way as the earlier example, you should be prepared to pay an infinite amount of money to play this game. This is obviously absurd. Very few people would be prepared to pay more than a few dollars for the chance to play the St. Petersburg game. So, something must be wrong with the concept of expected value. But what?



**Exercise 7.8** How much would you pay for a chance to play the St. Petersburg game? Explain your answer.



The St Petersburg game was first proposed by the eighteenth century Swiss mathematician Nicholas Bernoulli. His brother Daniel came up with a solution (which he published in the *Commentaries of the Imperial Academy of Sciences in St Petersburg* – hence the name). In essence, Daniel Bernoulli concluded that we should not try to value a gamble in terms of its expected value. In its place, he introduced the concept that we now call *utility* and proposed, in effect, that we think in terms of expected utility rather than expected value.

The basic idea behind the concept of utility is that utility is an index of the strength of your preference. In other words, to say that you assign more utility to X than to Y is to say that you prefer X to Y and so, given the choice, you will choose X over Y. The point of introducing the concept of utility is that, in the standard phrase, utility is not linear with money. That means that the utility you assign to something is not directly proportional to its expected (monetary) value.

Money has what is often called *diminishing marginal utility*. The additional utility you get from an extra \$10 depends upon how much money you already have. The \$10 that takes your net worth from \$100 to \$110 will probably mean much more to you than the \$10 that takes you from \$1,000,000 to \$1,000,010.

Different disciplines understand utility in different ways. From the perspective of most economists, for example, utility is solely a measure of preference as revealed by the choices people make. It is a purely operational notion, describing the choices that people make. On this way of thinking about utility, to say that a person assigns more utility to X than to Y is simply to say that, if they are consistent in certain ways, then they will choose X over Y. Utility is just a *description* of choice behavior.

For many psychologists and other social scientists, on the other hand, utility is not purely operational. It is a genuine psychological quantity that explains why people make the choices that they do. If I say that a person assigns more utility to X than to Y, then I am not just predicting that they will choose X over Y, I am explaining why they will make that choice. On this view, therefore, utility is an *explanation* of choice behavior. In cognitive science, utility is typically understood as an explanatory construct, rather than as a purely descriptive one.

Combining the concept of utility with the theory of probability discussed earlier gives the Bayesian approach to practical reasoning and decision-making. The central notion is the idea of expected utility. As its name suggests, expected utility is just like expected value, except that the concept of value is replaced by the concept of utility.

To calculate the expected utility of a possible action, you need to start by identifying the different possible outcomes that might result from that action. Then you assign a utility to each of those outcomes. But you need to take into account not just how much you value the different outcomes, but also how likely each of them is. So, you need to assign a probability to each outcome. Then for each outcome you multiply its utility by its probability. Adding together the results of this operation for each outcome gives you the expected utility for the action.

Expected utility is much more generally applicable than expected monetary value. It can be applied in cases where there is an easily identifiable expected value. You can assign

an expected utility to a lottery ticket. You might have a lottery ticket with an expected value of \$1 (it might be one of ten tickets in a lottery with a \$10 prize, for example). But you might assign that lottery ticket a lesser utility than you assign to having \$1 in cash. That is actually how economists define being risk-averse. Or you might assign the lottery ticket a greater utility than having \$1 in cash – which would be the definition of being risk-loving.

And you can also apply the concept of expected utility in nonmonetary contexts, where it is not obvious that there is any expected value. In fact, if you are a strict Bayesian you will think that utility and probability are all that you need. This is because Bayesians think that rational decision-makers will always act in a way that *maximizes expected utility*. In other words, the rational thing to do in any situation is always to choose the action that maximizes expected utility.



**Exercise 7.9** You are considering two possible actions – going for a swim in the river and going to watch a movie. The weather is doubtful and there is a 40 percent chance of rain. In general, you prefer swimming to movie-going, but not when it rains. Swimming offers 10 units of utility (standardly called *utils*) when it is not raining, but only 3 utils in the rain. Going to the movies will give you 6 units of utility, irrespective of the weather. As a Bayesian decision-maker, concerned only to maximize expected utility, which action should you choose?

This seems a long way from cognitive science, you might think. What has this all got to do with how we think about cognition? Well, in the next section we will see how this model can be applied in a very unexpected way – to shed light on the behavior of neurons in a region of the parietal cortex

## Case Study: Neurons That Code for Expected Utility

The experiments that we will be looking at in this section all rely on recording the activity of single neurons in monkeys. The technology for doing this was developed in the 1950s by Herbert Jasper at the Montreal Neurological Institute in Canada and Edward Evarts at the National Institutes of Health in the United States. Tiny electrodes are inserted into the monkey brain. This can be done while the animal is awake because there are no pain or touch receptors in the brain. The electrodes are small and sensitive enough to detect the firing rates of individual neurons.

Interestingly, some of the earliest experiments to use single-cell recording in awake monkeys studied the parietal cortex. Vernon Mountcastle and his research group of neurophysiologists developed influential experimental paradigms for studying how monkeys react to visual stimuli. They used fruit juice rewards to train their monkeys to make specific responses to visual cues. These experiments started a lengthy debate about what the parietal cortex actually does. Mountcastle took the view that the parietal cortex's job was to issue motor commands. The other side of the debate was taken up



by Michael Goldberg (then at the National Institutes of Health) who hypothesized that the parietal cortex had more to do with directing attention to highlight particular regions of the visual field.

The protracted debate between Mountcastle and Goldberg is fascinating but primarily concerns us as the background for a set of experiments carried out by Paul Glimcher (of New York University) and Michael Platt (from the University of Pennsylvania). They originally thought that they were developing experiments to settle the debate between Mountcastle and Goldberg. As it turned out, however, they ended up finding intriguing evidence of neurons that have significant Bayesian characteristics.

To understand the experiments, you need to know a little about how the muscles around the eye work. An important part of what they do is move the eyes in order to compensate for our own movements. This allows us to have (relatively) stable perceptions of the world around us. They also move the eyes so that the high-resolution part of the eye (the *fovea*) is focused on interesting and important things in the environment. These gaze alignment movements come in two varieties:

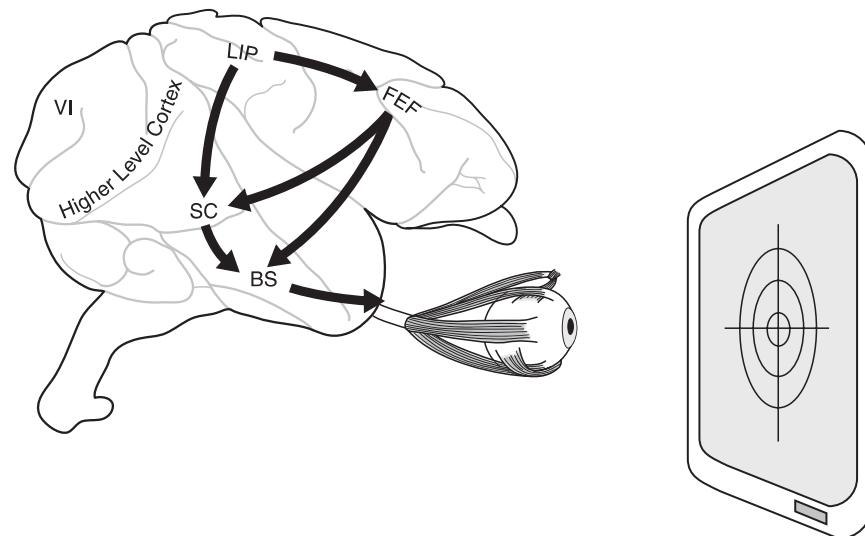
*Saccadic eye movements* move the line of sight very quickly from one place to another in the visual environment. This allows the perceiver to scan the environment (to detect a predator, for example)

*Smooth pursuit eye movements* allow the eyes to track objects moving continuously in a single direction.

Monkeys can be trained to perform saccadic eye movements. In a typical experiment, they are placed in front of a screen, fixating on a colored light directly in front of them. The experimenters flash a red spot on the right side of the screen. When the monkey made a saccadic eye movement toward the red spot, it is rewarded with a portion of fruit juice. The Glimcher and Platt experiments are basically variations on this basic theme.

What happens in the brain in between the monkey detecting a red spot on the right side of the visual field and its making a saccadic eye movement (a saccade) toward the red spot? Both ends of the process are relatively well understood. The primate visual system has been comprehensively mapped out, and so there is no mystery about how the monkey detects the red spot.

The other end of the process is how saccades are actually generated. This is also well understood. It is known that the *superior colliculus*, which is located in the midbrain, and the *frontal eye field*, which is in the frontal cortex, both play an important role in controlling saccades. It is also widely accepted that these two brain areas are organized topographically. That means that they are organized like a map, with individual neurons responsible for specific locations to which a saccade might be directed. So, just before the monkey makes a saccade to a specific location, the neuron corresponding to that location fires in the superior colliculus and/or frontal eye field.



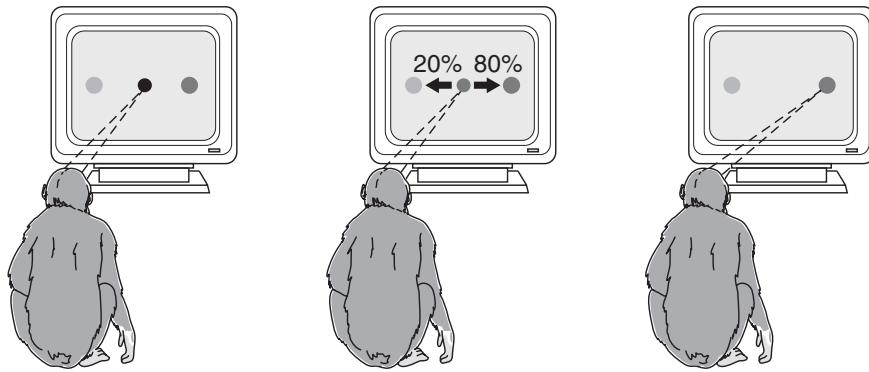
**Figure 7.6** The principal pathways for saccade production. LIP = lateral intraparietal area. FEF = frontal eye field. SC = superior colliculus. BS = brain stem. Note that, while LIP and the FEF are cortical structures, SC and BS are much more primitive areas and not part of the forebrain at all. (Figure 10.2 from Glimcher 2003: 230)

The key question, then is what happens in between detecting the red spot and initiating the saccade. This is where LIP (the lateral intraparietal area) comes in. Since LIP projects to both the superior colliculus and to the frontal eye field, it is the obvious place to look. Figure 7.6 illustrates the route of saccade production, showing why LIP is so important.

Platt and Glimcher came up with an ingenious set of experiments suggesting that LIP is essentially carrying out Bayesian calculations. The experiments were fairly complicated, but monkeys are extremely talented at learning things – and also highly motivated by fruit juice.

Experiments in this area are typically set up so that a particular action will always meet the same response. So, if a monkey is being trained to push a lever when it sees a red light on the right-hand side of its visual field, the correct action will always be rewarded (by delivery of fruit juice, typically). Moreover, the rewards are constant. If the reward for the correct action is 4 ml of juice on the first trial, it will still be 4 ml of juice on the twentieth trial.

Translating into our Bayesian language, what this means is that the reward is always delivered with probability 1, while the utility of the reward remains constant. Platt and Glimcher's breakthrough idea was to vary both probability and utility. Varying the size of the reward and how likely it is to be delivered allowed them to explore whether neurons are sensitive to those variations in the reward. In effect, it allowed them to test for Bayesian neurons.



**Figure 7.7** Platt and Glimcher's probabilistic cued saccade task. (Figure 10.11 from Glimcher 2003: 257)

## Probability-Detecting Neurons

To explore whether neurons in LIP are sensitive to probability, Platt and Glimcher set up a saccade experiment where the probability that the saccade would be rewarded varied.

As usual, the monkeys started by fixating on a light at the center of the screen. Shortly afterward, two other lights appeared on the screen, one on the left of the fixation point and one on the right. These were the targets for the saccades. Then the fixation light changed color, turning either red or green. The monkeys had been trained that when the fixation light turned red, they would be rewarded if they made a saccade to the left – while a saccade to the right would be rewarded if the light was green.

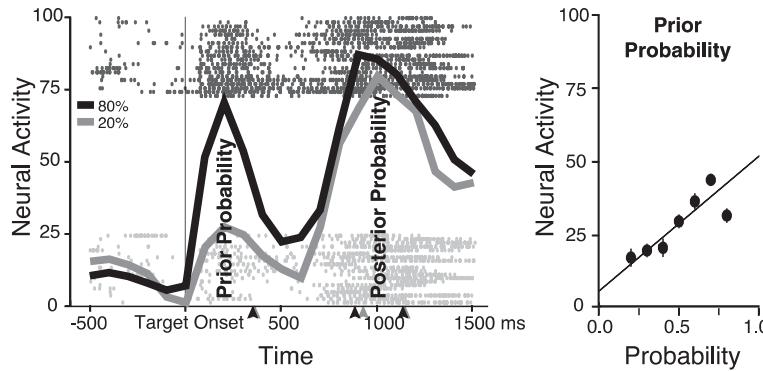
The twist to the experiment was that the fixation light was set up to turn red 80 percent of the time and green 20 percent (or the other way around, depending on the block). The experimental setup is illustrated in Figure 7.7.

Platt and Glimcher assumed that over the course of a 100-trial block the monkeys would have typically figured out that, say, a saccade to the left was much more likely to be rewarded than one to the right. So, they recorded throughout each trial from representative neurons in LIP.

To make sense of the results, they just compared, for each neuron, the trials where the stimulus and the movement were the same. So, for example, they compared all the trials in a given block where the monkey made a saccade to the left in response to a red light. Doing that made it possible to compare blocks where the probability of a red light (and hence a reward for making a saccade to the left) was low with blocks where red had a high probability.

Here's what they found (toward the end of the blocks, when the monkey should have learned which response was more likely to be rewarded).

- Prior to the fixation light changing color, in blocks where the probability of being rewarded with a left saccade was high, a typical LIP neuron had a much higher firing rate than in blocks where the probability of reward was low. Note that this is the firing rate *before* the monkey receives any indication as to which direction will be rewarded



**Figure 7.8** Activity of an LIP neuron during the probability experiment. Rows of tick marks in the panel indicate precise times of neural action potentials during each of twenty trials used to compute the averages shown as thick lines. (Figure 10.12 in Glimcher 2003: 260)

- After the fixation light changed color, the firing rate for the saccade that was definitely going to be rewarded was maximal, even if the probability that the saccade would be rewarded had been low. Note that this is the firing rate *after* the fixation light has indicated which direction will be rewarded.

Translating this back into Bayesian terms, the first finding seems to show that the neurons are encoding the prior probabilities – the probability that a given saccade will be rewarded, before it is known which saccade will be rewarded. And the second seems to show that the neurons are updating the priors to posteriors when the fixation light changes color. Once the fixation light has turned either red or green it is certain which saccade will be rewarded, no matter whether it was a high-probability saccade or a low-probability one. And the neurons respond by firing at full blast, as it were. Figure 7.8 illustrates both these findings.

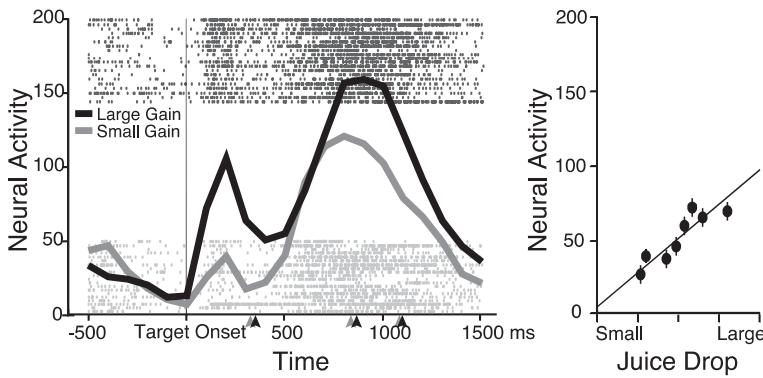
You can see how the neuron's firing rates in the low-probability and high-probability conditions are very different before the moment when the fixation point changes color. Then they subsequently end up firing at the same rates, because the changing color of the fixation point indicates that they actually will be rewarded (remember that all the trials illustrated ended with rewards).



**Exercise 7.10** Explain in your own words why Platt and Glimcher think that they have discovered evidence of neurons that code for probability.

## Utility-Detecting Neurons

So – what happens if probabilities (and all the other dimensions of the experiment) are held constant, but the quantity of the reward is varied? Well, as Platt and Glimcher observed, that would be a very good test of whether LIP neurons are sensitive to utility.



**Figure 7.9** Platt and Glimcher's cued saccade experiment, with stimulus and response held constant and the quantity of reward varied. The neuron's average firing rate is shown for a high-reward block (bold line) and a low-reward block (light-colored line). (Figure 10.13 from Glimcher 2003: 262)

And to explore it they used a very similar cued saccade task to the one we have been considering.

As before, the monkeys fixate on a point in the center of the screen. When the fixation point changes color, that indicates to the monkey that the reward will be delivered if the monkey makes a saccade to the left (red) or the right (green). The probability of red versus green was held constant, with each color coming up on exactly half the trials. What varied was the quantity of the reward. In one block, for example, the reward for looking left might be double that for looking right.

The results are illustrated in Figure 7.9. As before, what the figure shows are the average response profiles over time of a single neuron in two different blocks. In both blocks the stimulus and response are held constant. They are all cases, say, where the light turns green and the monkey is rewarded for a saccade to the right. What varies is the quantity of the reward. In one block, marked by a light-colored line, the response receives a low reward. In the other, marked with a bold line, the response receives a high reward.

Before the light changes color, the neuron fires more strongly on average in the high-reward block than in the low-reward block. After the light changes color (and so when the reward is revealed), the average firing rate increases significantly in both conditions. But the difference across the two conditions remains constant. The neuron responds more vigorously to the larger reward.

Strictly speaking (as Platt and Glimcher note), this experiment does not show that LIP neurons are sensitive to utility. Their firing rates correlate with the quantity of reward, but that is not necessarily the same as the utility that the monkey might derive from the reward. The experiments don't reveal any analog to the phenomenon of diminishing marginal utility, for example. But still, they are highly suggestive, and exactly what one would expect in a Bayesian brain!

## Combining Probability and Utility

So, the two different cued saccade experiments just reviewed seem to uncover neurons in LIP that code, first, for the probability of a reward and, second, for the size of a reward (and perhaps for its utility). We are close to having all the ingredients for a fully Bayesian brain, in which neurons code for something very close to expected utility. Because, as we saw earlier, expected utility is really just a combination of probability and utility. The expected utility of an action is arrived at by summing the utility of its different possible outcomes, each weighted by the probability with which it will occur. That takes us to the final experiments for this case study, because Platt and Glimcher came up with intriguing results here too – not quite expected utility, but intriguingly close.

To study sensitivity to expected utility, Platt and Glimcher needed to adapt the cued saccade paradigm we have been looking at. They needed to introduce an element of choice, because expected utility really only applies when monkeys (or people) have more than one available action. The principle of expected utility is a tool for choosing between two or more available actions (choose the action with the highest expected utility).

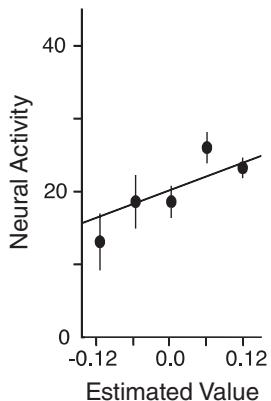
So, they turned the cued saccade task into a free-choice task. As before, the monkey fixated on a spot at the center of the screen. But this time he could choose to look at a stimulus on the left or a stimulus on the right. Everything was held constant within each block of 100 trials, but the quantity of the (fruit juice) reward varied across blocks – 0.2 ml for looking left and 0.1 ml for looking right, for example.

From a strict Bayesian perspective, there is an optimal way of responding in this type of experiment. You should sample the two alternatives until you have worked out which yields the highest reward, and then stick with that option until the end of the trial. That would be the best way to maximize expected utility over the long run.

It turns out, though, that the monkeys did not adopt the optimal strategy. Instead they displayed a form of *matching behavior*. That is, they split their choices between the two alternatives in a way that matched the distribution of total reward across the two alternatives. In other words, if looking right yielded twice the reward of looking left, then they looked right twice as often as they looked left. This is not the behavior that you would expect from an animal that was maximizing expected utility. However, things turned out to be more complicated than initially appeared.

Platt and Glimcher noticed something very interesting when they looked at what the monkeys were doing on a trial-by-trial basis. They saw that there was a pattern to the individual “choices” that the monkeys made. When they looked at each individual choice, they saw that the monkeys seemed to be engaged in a form of maximizing behavior. The monkeys were not maximizing expected utility, but they seemed to be maximizing something else that was not that far off from expected utility.

They used a version of *melioration theory*, developed by the animal behaviorist Richard Herrnstein (who, incidentally, first discovered the matching phenomenon when studying pigeon pecking behavior) in order to estimate how much value the monkey attached to



**Figure 7.10** Activity of an LIP neuron while a monkey makes his own choice compared to a behaviorally derived estimate of the value of the movement to the monkey. (Figure 10.15 from Glimcher 2003: 266)

each option in each trial. So, they were able to estimate the expected value to the monkey of looking left versus looking right on trial 47, for example.

According to melioration theory, value is fixed by local rates of reward. An animal behaving according to melioration theory will choose the option that seems most attractive at that time. And attractiveness at a time is fixed by the rewards that that option has yielded in recent trials. So, if looking left has had a good track record in yielding the reward on the past few trials, then the monkey will continue to look left. But if looking left starts to look unsuccessful, the monkey will switch and start to look right.

It turned out that the monkeys were actually behaving consistently with melioration theory, when each trial was viewed in the context of the previous ten trials. What that means is that the monkeys tended to choose the option that had been most highly rewarded over the previous ten trials. So, they were maximizing expected value, as calculated by melioration theory. This is not the same as the expected utility of each of the two options, because it is a local quantity that looks back to the history of rewards, whereas expected utility is more of a global measure of overall rewards. But still, it is certainly a related measure (and in fact, expected value in the melioration theory sense, is a good approximator of expected utility in many contexts). This is illustrated in Figure 7.10.

The important point from the perspective of neural economics, though, is that when Platt and Glimcher looked at the data from this perspective, they discovered a close correlation between the firing rates of individual neurons and the anticipated value to the animal, as computed by melioration theory. So, when the estimated value to the animal of a particular option was low (based on the results of the previous ten trials) the firing rates would be low. And the greater the estimated value, the higher the firing rate. The neurons, it seemed, were firing in accordance with expected value, as calculated by melioration theory.

In conclusion, the experiments do not show that neurons in LIP are calculators of expected utility. But they do seem to uncover neurons sensitive to the varying benefits to the animal of different courses of action. Platt and Glimcher can have the last word:

In our free-choice task, both monkeys and posterior parietal neurons behaved as if they had knowledge of the gains associated with different actions. These findings support the hypothesis that the variables that have been identified by economists, psychologists and ecologists as important in decision-making are represented in the nervous system.

(Platt and Glimcher 1999)



## Summary

This chapter began by introducing the basic tools for understanding Bayesian approaches to cognitive science. Reviewing the elements of the probability calculus allowed us to introduce Bayes's Rule as a tool for measuring the strength of support that evidence provides for a hypothesis. Our first example of how the Bayesian tool kit can be applied within cognitive science came from visual perception. The perceptual systems can be modeled as continually solving a Bayesian problem. They have to arrive at a hypothesis about the layout of the distal environment on the basis of partial and noisy sensory data (the evidence), and so Bayes's Rule seems an ideal tool. We illustrated this approach by looking at a Bayesian model of binocular rivalry.

The second part of the chapter explored Bayesian approaches to choice behavior and practical decision-making. The key theoretical tool here is the concept of utility, which measures how much a decision-maker values a particular outcome. Bayesian theories of choice typically rely upon the principle of expected utility, which states that rational decision-makers will choose actions that maximize expected utility, where the expected utility of an action is given by the utility of its possible outcomes, each weighted by the probability of that outcome. We looked at a series of experiments in the burgeoning field of neuroeconomics, illustrating how neural activity can be modeled in Bayesian terms. Platt and Glimcher's recordings of single neurons in the LIP area of the monkey parietal cortex uncover neurons that are sensitive to probability and to analogs of utility and expected utility.

## Checklist

### **Bayesianism is built on three basic ideas:**

- (1) Belief comes in degrees
- (2) Degrees of belief can be modeled as probabilities, which means that they have to conform to the basic principles of the probability calculus



- (3) Learning takes place by updating probabilities according to Bayes's Rule.

**Bayes's Rule is a tool for measuring the posterior probability of a hypothesis, conditional upon some evidence:**

- (1) Bayes's Rule is calculated from the *prior probability* of the hypothesis,  $p(H)$ , the *likelihood* of the evidence, given the hypothesis,  $p(E/H)$ , and the probability of the evidence,  $p(E)$ .
- (2) According to Bayes's Rule, the posterior probability of the hypothesis =

$$\frac{\text{Likelihood of the evidence} \times \text{Prior probability of the hypothesis}}{\text{Probability of the evidence}}$$

- (3) Bayes's Rule is a straightforward consequence of the definition of conditional probability.

**Perception can be modeled in terms of Bayesian inference.**

- (1) Perceptual systems have the job of selecting between different hypotheses about the layout of the distal environment on the basis of noisy and partial data (the evidence).
- (2) This process can be modeled as an application of Bayes's Rule, because ultimately perceptual systems have to decide which hypothesis has the highest posterior probability, conditional upon the evidence coming from the senses.

**Binocular rivalry offers a case study in Bayesian approaches to perception.**

- (1) When separate images are presented to each eye, the percept generated by the visual system alternates between the two images, seemingly at random.
- (2) From a Bayesian perspective, this is an understandable response to the posterior probabilities calculated via Bayes's Rule.
- (3) The posterior probability for hypothesis corresponding to the left image (conditional upon the evidence) is the same as the posterior probability for the right image hypothesis, and both are higher than the posterior probability for the hypothesis that the environment contains a composite corresponding to a blend of the two images.
- (4) Since the posterior probabilities for the left and right hypotheses are equal, the visual system can't decide between them, and so simply alternates.

**Expected utility is the key concept for applying Bayesianism to practical decision-making.**

- (1) Utility is a measure of how much a decision-maker values a particular outcome.
- (2) To calculate the expected utility of an action you need to assign utilities to its different possible outcomes and then those utilities together, each discounted by the probability that it will occur.
- (3) The St Petersburg game shows that expected utility is not necessarily the same as expected monetary value.

**Platt and Glimcher's experiments recording single neurons in the parietal cortex of the monkey are an illustration from neuroeconomics of how the brain may apply Bayesian principles.**

- (1) Neuroeconomics is located at the interface between neuroscience, microeconomics, the psychology of reasoning, behavioral finance, and decision theory.
- (2) Platt and Glimcher studied how individual neurons in LIP responded while monkeys made saccadic eye movements for fruit juice rewards.
- (3) Holding all other aspects of the task constant, but varying the probability that each response would be rewarded revealed that the firing rates of individual neurons correlates with both the prior and posterior probabilities of reward.
- (4) Holding all aspects of the task constant, but varying the quantity of the reward revealed that the firing rates of individual neurons correlates with the value of the reward (an analog of utility)
- (5) Using a free-choice version of the cued saccade task revealed neurons firing according to the estimated value to the animal of a particular response (as calculated according to Herrnstein's melioration theory, which differs from but is related to expected utility theory).

## Further Reading

Good general introductions to Bayesian approaches to cognitive science can be found in articles by Griffiths, Kemp, and Tenenbaum (2008), Chater et al. 2010, and Jacobs and Kruschke 2011. For an overview of specific applications, see the special issue of *Trends in Cognitive Sciences* from 2006, devoted to "Probabilistic models of cognition." The essays in Chater and Oaksford 2008 cover a wide range of Bayesian models.

Skyrms 1986 is a classic introduction to inductive logic and theories of probability. Hacking 2001 is more user-friendly. See Kaplan 1996 for a general philosophical defense of Bayesian approaches to belief and Horwich 1982 for a Bayesian approach to the philosophy of science. Jeffrey 1983 is a very influential presentation of Bayesian decision theory. For a historical introduction to different ways of thinking about utility, see Broome 1991. I have written about Bayesian decision theory as a theory of rationality in Bermúdez 2009.

For Bayesian approaches to perception in general, see the papers in Knill and Whitman 2008, and for an article-length review focused on visual perception, see Kersten, Mamassian, and Yuille 2004. Hohwy 2013 uses visual perception, and in particular, the example of binocular rivalry discussed in Section 7.2 as the starting point for a general theory of the mind as a Bayesian predictive machine. The Bayesian model of binocular rivalry is presented in more detail in Hohwy, Roepstorff, and Friston 2008. Clark 2016 discusses predictive coding, relating it to embodied and situated cognition. For further good illustrations of Bayesian approaches, see Wolpert's work on motor control (Kording and Wolpert 2004, 2006) and Ernst and Banks on multisensory integration (Ernst and Banks 2002).



There is an overview of neuroeconomics from a philosophical perspective in Hardy-Vallé 2007 and articles on a range of applications in a special issue of *Brain Research Bulletin* from November 2005. The principal textbook for the field is Glimcher and Fehr 2014 (2nd ed.). Mountcastle et al. 1975 and Robinson, Goldberg, and Stanton 1978 are important early papers in the debate about the specific role of neurons in the posterior parietal cortex. Chapter 10 of Glimcher 2003 accessibly tells the story of single-neuron studies of LIP, leading up to his own experiments, which are presented more rigorously in Platt and Glimcher 1999.





## CHAPTER EIGHT

# Modules and Architectures

### OVERVIEW 203

#### 8.1 Architectures for Artificial Agents 204

Three Agent Architectures 204

#### 8.2 Fodor on the Modularity of Mind 208

Modular and Nonmodular Processing 208

#### 8.3 The Massive Modularity Hypothesis 210

The Cheater Detection Module 211

The Evolution of Cooperation 213

Two Arguments 216

Evaluating the Arguments for Massive Modularity 218

#### 8.4 Hybrid Architectures: The Example of ACT-R 219

The ACT-R Architecture 220

ACT-R as a Hybrid Architecture 222



## Overview

This chapter tackles the overall organization of the mind. We start in Section 8.1 by looking at *agent architectures* in AI. These are blueprints for the design of artificial agents. Artificial agents can be anything from robots to internet bots. Looking at different architectures allows us to distinguish cognitive systems from, for example, reflex systems, or reflex agents. Reflex systems are governed by simple production rules that uniquely determine how the system will behave in a given situation. In contrast, cognitive systems deploy information processing between the input (sensory) systems and the output (effector) systems.

Intelligent agents in AI are standardly built up from subsystems that perform specific information-processing tasks. Cognitive scientists tend to think of the mind (at least in part) as an organized collection of specialized subsystems carrying out specific information-processing tasks. The earliest sustained development of this idea from a theoretical point of view came in a book entitled *The Modularity of Mind*, written by the philosopher Jerry Fodor. We look at Fodor's modularity thesis in Section 8.2.

Fodor distinguishes modular processing from central processing, responsible for general problem solving and decision-making. Massive modularity theorists, in contrast, deny that there is any such thing as nonmodular central processing. We look at this model in Section 8.3 and see how it has been used to explain research into the psychology of reasoning.



Finally, in Section 8.4 we relate the discussion of agent architectures back to earlier discussions of different ways of thinking about information and information processing. We look at an example of a hybrid architecture combining the physical symbol system hypothesis and the parallel processing characteristic of connectionist networks. This is the ACT-R architecture, developed by John R. Anderson and colleagues at Carnegie Mellon University.



## 8.1

# Architectures for Artificial Agents

One aim of AI researchers is to build artificial agents. There are many different types of AI agents. Robots are probably the first things to come to mind when thinking about intelligent agents. Robotic agents are built to operate in real, physical environments. But many agents are designed to function in virtual environments. Shopping bots are good examples. Some bots are designed to travel around the internet comparing prices for a single item, while others trawl through sites such as Amazon finding items that you might be likely to buy (perhaps because they have been bought by customers who bought some items that you bought).

Computer scientists have come up with an interesting range of different *agent* architecture for designing artificial agents. An agent architecture is a blueprint that shows the different components that make up an agent and how those components are organized. In this section we will look at three different types of agent architecture:

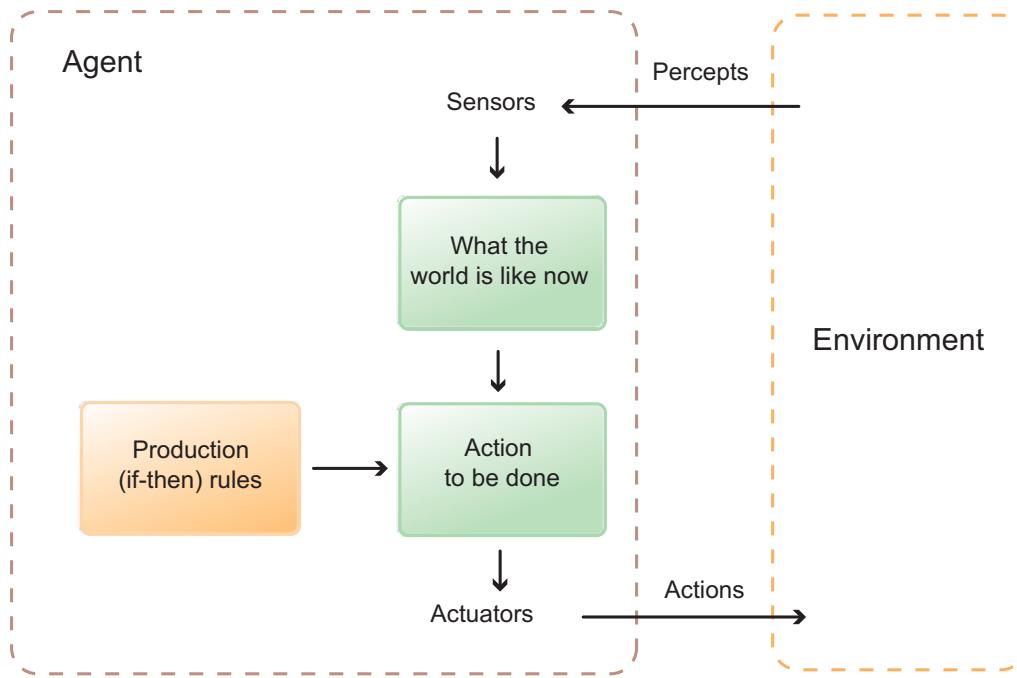
- A simple reflex agent
- A goal-based agent
- A learning agent

Not all artificial agents are intelligent agents. Looking at these architectures reveals what is distinctive about intelligent, cognitive agents, as opposed to simpler, noncognitive agents. The agent architectures we will be looking at range from the plainly noncognitive to the plainly cognitive. As we go through them we get a better picture of the basic functions that any cognitive system has to perform.

First, we need to know what an agent is. The quick definition is that an agent is a system that perceives its environment through *sensory systems* of some type and acts upon that environment through *effector systems*. The basic challenge for a computer scientist programming an agent (whether a software agent or a robotic agent) is to make sure that what the agent does is a function of what the agent perceives. There need to be links between the agent's sensory systems and its effector systems. What distinguishes different types of agent is the complexity of those links between sensory systems and effector systems.

## Three Agent Architectures

The simplest type of agent in agent-based computing is the *reflex agent*. Simple reflex agents have direct links between sensory and effector systems. The outputs of the sensory systems



**Figure 8.1** The architecture of a simple reflex agent. Production rules are all that intervene between sensory input and motor output. (Adapted from Russell and Norvig 2009)

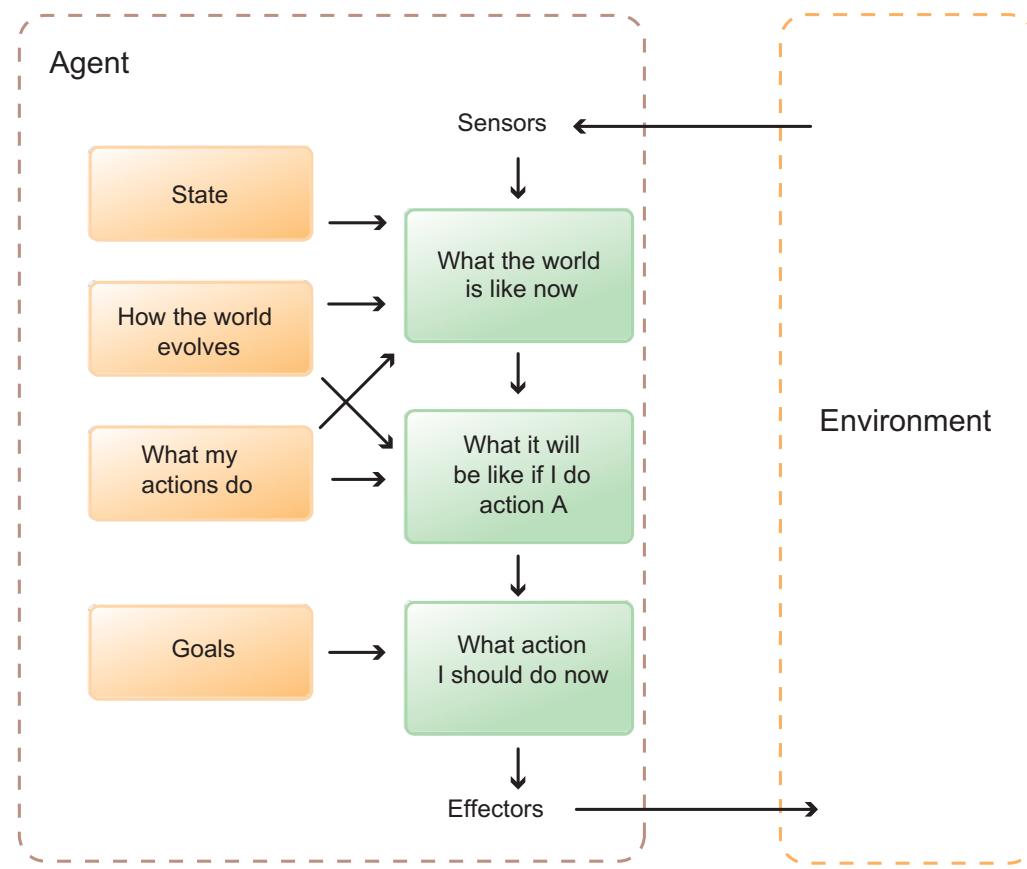
directly determine the inputs to the effector systems. These direct links are achieved by *production rules*. Production rules take this form:

IF condition C holds THEN perform action A.

It is up to the sensory systems to determine whether or not condition C holds. Once the sensory systems have determined that condition C holds, then the behavior of the simple reflex agent is fixed. Figure 8.1 shows a schematic representation of the architecture of a simple reflex agent.

Simple reflex agents are not, many cognitive scientists would think, cognitive systems. This is because they are simply reacting to the environment in invariant ways – the same stimulus always receives the same response. In contrast, it is often taken to be an essential feature of cognitive systems that they can react differently to the same environmental stimulus. This is because the actions of cognitive systems are determined by their goals and by their stored representations of the environment. Human agents, for example, sometimes act in a purely reflex manner. But more often we act as a function of our beliefs and desires – not to mention our hopes, fears, dislikes, and so on.

The schematic agent architecture in Figure 8.2 depicts a primitive type of cognitive system. This is a *goal-based agent*. As the diagram shows, goal-based agents do not simply act upon environmental stimuli. There are no simple production rules that will uniquely determine how the agent will behave in a given situation. Instead, goal-based agents need to work out the consequences of different possible actions and then evaluate those



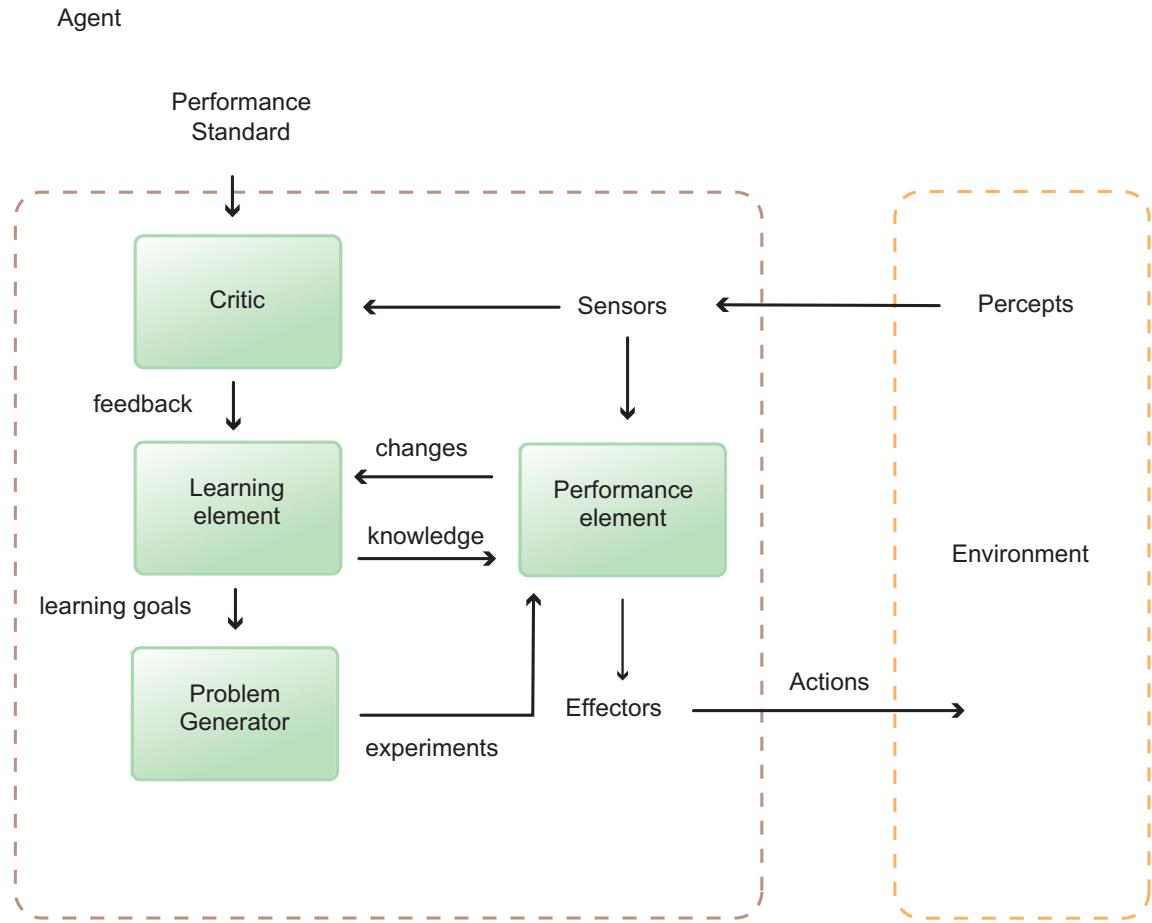
**Figure 8.2** The architecture of a goal-based agent. There are information-processing systems intervening between input and output. (Adapted from Russell and Norvig 2009)

consequences in the light of their goals. This is done by the specialized cognitive systems labeled in Figure 8.2.

But there is still something missing from goal-based agents. They have no capacity to learn from experience – something which, surely, is necessary for any agent to count as an intelligent agent. For that we need a *learning agent*. A sample architecture for a *learning agent* is presented in Figure 8.3.

The learning agent has certain standards that it wants its actions to meet. These are one of the inputs to the Critic subsystem, which also receives inputs from the sensory systems. The Critic's job is to detect mismatches between sensory feedback and the performance standard. These mismatches feed into the Learning subsystem which determines learning goals and makes it possible for the system to experiment with different ways of achieving its goals.

As the learning agent example shows, computer scientists designing intelligent agents typically build those agents up from subsystems performing specific information-processing tasks. This way of thinking about cognitive systems (as organized complexes



**Figure 8.3** The architecture of a learning agent. Mismatches between sensory feedback and the performance standards are detected by the Critic subsystem. The Learning subsystem determines learning goals and allows the system to experiment with different ways of achieving its goals. (Adapted from Russell and Norvig 2003)

of subsystems) has proved very influential in cognitive science. It raises a number of important questions:

- How are we to identify and distinguish cognitive subsystems?
- Are there any important differences between the subsystems responsible for sensory processing and motor behavior, on the one hand, and those that operate between those input and output subsystems?
- Do all the subsystems in a cognitive system process information in the same way? Do they all involve the same type of representations?
- How “autonomous” are the different subsystems? How “insulated” are they from each other?

To explore these questions further we turn now to the hypothesis that the mind is modular, proposed by the influential philosopher and cognitive scientist Jerry Fodor in his book, *The Modularity of Mind*, published in 1983.


**8.2**

## Fodor on the Modularity of Mind

### Modular and Nonmodular Processing

Fodor proposes a radical distinction between modular and nonmodular cognitive processes. On the one hand, nonmodular processes are high-level, open-ended, and involve bringing a wide range of information to bear on very general problems. In contrast, modular processes operate at a much lower level and work quickly to provide rapid solutions to highly determinate problems.

As Fodor defines them, modular processes have the following four characteristics:

- *Domain-specificity.* Modules are highly specialized. They are mechanisms designed to carry out very specific and circumscribed information-processing tasks. Because of this they only operate on a limited range of inputs (those relevant to their particular domain)
- *Informational encapsulation.* Modular processing is not affected by what is going on elsewhere in the mind. Modules cannot be “infiltrated” by background knowledge and expectations, or by information in the databases associated with different modules.
- *Mandatory application.* Cognitive modules respond automatically to stimuli. They are not under any executive control and cannot be “switched off.” It is evidence that certain types of visual processing are modular that we cannot avoid perceiving visual illusions, even when we know them to be illusions.
- *Speed.* Modular processing transforms input (e.g., patterns of intensity values picked up by photoreceptors in the retina) into output (e.g., representations of three-dimensional objects) quickly and efficiently.

In addition to these “canonical” characteristics of modular processes, Fodor draws attention to two further features that sometimes characterize modular processes.

- *Fixed neural architecture.* It is sometimes possible to identify determinate regions of the brain associated with particular types of modular processing. So, for example, an area in the fusiform gyrus (the so-called fusiform face area) is believed to be specialized for face recognition, which is often described as a modular process.
- *Specific breakdown patterns.* Modular processing can fail in highly determinate ways. These breakdowns can provide clues as to the form and structure of that processing. Prosopagnosia is a highly specific neuropsychological disorder that affects face recognition abilities, but not object recognition more generally.

These last two features are less central, because cognitive scientists tend to think of cognitive modules in terms of their function (the information-processing task that they carry out), rather than their physiology. A cognitive module has to perform a single, circumscribed, domain-specific task. But it is not necessary that it map onto a particular



part of the brain. Some modules do seem to be localizable, but for others we have (as yet) no evidence either way.

Cognitive modules form the first layer of cognitive processing. They are closely tied to perceptual systems. Here are some mechanisms that Fodor thinks are likely candidates for cognitive modules:

- Color perception
- Shape analysis
- Analysis of three-dimensional spatial relations
- Visual guidance of bodily motions
- Grammatical analysis of heard utterances
- Detecting melodic or rhythmic structure of acoustic arrays
- Recognizing the voices of conspecifics

Some of these candidate modules are close to the sensory periphery. In other words, relatively little information processing occurs between the sense organs and the module. This is clearly the case for color perception. Other systems are much further “downstream.” An example here would be grammatical analysis. A lot of processing needs to take place before there is an auditory or visual image to be analyzed for grammatical structure. Moreover, some cognitive modules can take the outputs of other modules as inputs. It is likely that information about the rhythmic structure of an acoustic array will be relevant to identifying the voice of a conspecific, for example.

Most of the modules Fodor discusses are involved in perceptual information processing, but it seems likely that many motor tasks are also carried out by modules. Planning even the simplest reaching movement involves calibrating information about a target object (a glass, say) with information about hand position and body orientation. The location of the glass needs to be coded on a hand-centered coordinate system (as opposed to one centered on the eyes, for example). Executing the movement requires, first, calculating a trajectory that leads from the start location to the end location, and then calculating an appropriate combination of muscle forces and joint angles that will take the arm along the required trajectory. These are all highly specialized tasks that seem not to depend upon background information or central processing – prime candidates for modular processing, on Fodor’s analysis.

But, according to Fodor, not all cognition can be carried out by modular mechanisms. He insists that there have to be psychological processes cutting across cognitive domains. The very features of cognitive modules that make them computationally powerful, such as their speed and informational encapsulation, mean that their outputs are not always a good guide to the layout of the perceived environment. Appearances can be deceptive. This means that there has to be information processing that can evaluate and correct the outputs of cognitive modules. This is what Fodor and others often call central processing, to distinguish it from modular processing, which is peripheral.

For Fodor, central processing has two distinguishing features. The first distinguishing feature is that central processing is *Quinean* (named after the philosopher Willard von Orman Quine, who famously proposed a holistic view of knowledge and confirmation). For Fodor, each organism's belief system is, in important respects, analogous to a scientific theory. It is, in fact, the organism's theory of the world, and so shares important properties with scientific theories. It is the belief system as a whole that is evaluated for consistency and coherence, for example. We cannot consider how accurate or well confirmed individual beliefs are in isolation, since how we evaluate individual beliefs cannot be divorced from how we think about other elements of the system in which they are embedded.

The second distinguishing feature of central processing is that it is *isotropic*. The isotropic nature of central processing is in many ways a corollary of its Quinean property. To say that central processing is isotropic is, in essence, to say that it is *not* informationally encapsulated. In principle any part of the belief system is relevant to confirming (or disconfirming) any other. We cannot draw boundaries within the belief system and hope to contain the process of (dis)confirmation within those boundaries.

Fodor himself is very pessimistic about cognitive science being able to shed any light on central processing. Cognitive science, Fodor argues, is really best suited to understanding modular processes. It can tell us very little about central processes – about all the processing that takes place in between sensory systems and motor systems. Unsurprisingly, this is not a view that has gained much currency within cognitive science as a whole.

Nonetheless, there are influential movements within cognitive science that are skeptical of the importance and even the existence of what Fodor calls central processing. We have already seen some examples of this from dynamical systems theory in Chapter 6. The next section looks at a very different approach – the massive modularity hypothesis.

### 8.3 The Massive Modularity Hypothesis

Supporters of the massive modularity hypothesis claim that the mind does not really do any central processing at all. They think that all information processing is essentially modular, although they understand modules in a much less strict way than Fodor does.

According to the massive modularity hypothesis, the human mind is a collection of specialized modules, each of which evolved to solve a very specific set of problems that were confronted by our early ancestors – by hunter-gatherers in the Pleistocene period, or even earlier in the evolutionary history of the human ape. These are called *Darwinian modules*.

What sort of Darwinian modules might there be? Evolutionary psychologists have tended to focus primarily on modules solving problems of social coordination, such as problems of cheater detection, kin detection, and mate selection. But massive modularity theorists are also able to appeal to evidence from many different areas of cognitive science pointing to the existence of specialized cognitive systems for a range of different abilities and functions. These include:



- Face recognition
- Emotion detection
- Gaze following
- Folk psychology
- Intuitive mechanics (folk physics)
- Folk biology

There is some overlap between Fodorean modules and Darwinian modules, but typically Darwinian modules engage in more complex types of information-processing than Fodorean ones. As a consequence, they are typically not informationally encapsulated. Massive modularity theorists describe Darwinian modules as modular primarily because they are domain-specific.

Many different types of evidence are potentially relevant to identifying Darwinian modules. In Chapter 11 we will look at some influential experiments on prelinguistic infants using the dishabituation paradigm. These experiments show that infants are perceptually sensitive to a number of basic principles governing the behavior of physical objects – such as the principle that objects follow a single continuous path through space and time. These experiments have been taken to show that infants possess a basic theory of the physical world. This basic theory is held by many to be the core of adult folk physics, which itself is the domain of a specialized cognitive system.

The case for massive modularity rests upon a mixture of case studies and general arguments. The most developed case study is the so-called cheater detection module. We will look at this first, as a detailed illustration of a Darwinian module. Then we will turn to the general arguments.

## The Cheater Detection Module

We need to start from well-known experiments on reasoning with conditionals (sentences that have an IF ... THEN ... structure). These experiments, often using variants of a famous experiment known as the Wason selection task, have been widely taken to show that humans are basically very poor at elementary logical reasoning. It turns out, however, that performance on these tasks improves drastically when they are reinterpreted to involve a particular type of conditional. These are so-called *deontic conditionals*. Deontic conditionals have to do with permissions, requests, entitlements, and so on. An example of a deontic conditional would be: If you are drinking beer then you must be over 21 years of age.

The evolutionary psychologists Leda Cosmides and John Tooby came up with a striking and imaginative explanation for the fact that humans tend to be much better at reasoning with deontic conditionals than they are with ordinary, nondeontic conditionals. According to Cosmides and Tooby, when people solve problems with deontic conditionals they are using a specialized module for monitoring social exchanges and detecting cheaters. This is the cheater detection module.

E C 4 5

**Figure 8.4** A version of the Wason selection task. Subjects are asked which cards they would have to turn over in order to determine whether the following conditional is true or false: **If a card has a vowel on one side, then it has an even number on the other.**

To see how this works, let's start with a typical version of the Wason selection task. Imagine that you are shown the four cards illustrated in Figure 8.4 and told that each card has a letter on one side and a number on the other. The experimenter then asks you which cards you would need to turn over in order to determine whether the following conditional is true or false: **If a card has a vowel on one side then it has an even number on the other.**

It is obvious that the *E* card will have to be turned over. Since the card has a vowel on one side, the conditional will certainly be false if it has an odd number on the other side. Most subjects get this correct. It is fairly clear that the second card does not need to be turned over, and relatively few subjects think that it does need to be turned over. The problems arise with the two numbered cards.

Reflection shows (or should show!) that the *4* card does not need to be turned over, because the conditional would not be disconfirmed by finding a consonant on the other side. The conditional says that any card with a vowel on one side has to have an even number on the other side. It doesn't say anything about cards that have consonants on one side and so the existence of a card with a consonant on one side and an even number on the other is irrelevant.

The *5* card, however, does need to be turned over, because the conditional will have to be rejected if it has a vowel on the other side (this would be a situation in which we have a card with a vowel on one side, but no even number on the other).

Unfortunately, very few people see that the *5* card needs to be turned over, while the vast majority of experimental subjects think that the *4* card needs to be turned over. This result is pretty robust, as you will find out if you try it on friends and family.

So, what is going wrong here? It could be that the experimental subjects, and indeed the rest of us more generally, are reasoning in perfectly domain-general ways, but simply employing the wrong domain-general inferential rules. But further work on the Wason suggestion task has suggested that this may not be the right way of thinking about it.

It turns out that performance on the selection task varies drastically according to how the task is formulated. There are "real-world" ways of framing the selection task on which the degree of error is drastically diminished. One striking set of results emerged from a variant of the selection task carried out by Richard Griggs and James Cox. They transformed the selection task from what many would describe as a formal test of conditional reasoning to a problem-solving task of a sort familiar to most of the experimental subjects.

Griggs and Cox preserved the abstract structure of the selection task, asking subjects which cards would have to be turned over in order to verify a conditional. But the conditional was a



**Figure 8.5** A version of Griggs and Cox's deontic selection task. Subjects are asked to imagine that they are police officers checking for underage drinkers and asked which cards they would need to turn over in order to assess the following conditional: **If a person is drinking beer, then that person must be over 19 years of age.**

conditional about drinking age, rather than about vowels and even numbers. Subjects were asked to evaluate the conditional: **If a person is drinking beer, then that person must be over 19 years of age** (which was, apparently, the law at the time in Florida). They were presented with the cards shown in Figure 8.5 and told that the cards show the names of drinks on one side and ages on the other. Before making their choice, subjects were told to imagine that they were police officers checking whether any illegal drinking was going on in a bar.

The correct answers (as in the standard version of the selection task we have already considered) are that the *BEER* card and the *16* card need to be turned over. On this version of the selection task, subjects overwhelmingly came up with the correct answers, and relatively few suggested that the third card would need to be turned over. What is particularly interesting is the subsequent discovery that if the story about the police officers is omitted, performance reverts to a level comparable to that on the original selection task.

The finding that performance on the selection task can be improved by framing the task in such a way that what is being checked is a condition that has to do with permissions, entitlements, and/or prohibitions has proved very robust. The fact that we are better at reasoning with these *deontic* conditionals than we are with ordinary conditionals has suggested to many theorists that we have a *domain-specific* competence for reasoning involving permissions and prohibitions.

Building on these results, the evolutionary psychologists Leda Cosmides and John Tooby have suggested that the human mind (perhaps in common with the minds of other higher apes) has a dedicated cognitive system (a *module*) for the detection of cheaters. The cheater detection module is supposed to explain the experimental data on the Wason selection task. When the selection task is framed in terms of permissions and entitlements it engages the cheater detection module. This is why performance suddenly improves.

## The Evolution of Cooperation

But why should there be a cheater detection module? What was the pressing evolutionary need to which the cheater detection module was a response? Cosmides and Tooby's answer these questions through an influential theory of the emergence of cooperative behavior.

Biologists, and evolutionary theorists more generally, have long been puzzled by the problem of how cooperative behavior might have emerged from a process of natural

selection. Cooperative behavior presumably has a genetic basis. But how could the genes that code for cooperative behavior ever have become established, if (as seems highly plausible) an individual who takes advantage of cooperators without reciprocating will always do better than one who cooperates? Evolution seems to favor free riders and exploiters above high-minded altruists.

A popular way of thinking about the evolution of cooperation is through the model of the prisoner's dilemma. The prisoner's dilemma is explained in Box 8.1. Many interpersonal interactions (and for that matter many interactions between nonhuman animals) involve a series of encounters each of which has the structure of a prisoner's dilemma, but where it is not known how many encounters there will be. Game theorists call these indefinitely iterated prisoner's dilemmas.

One way of dealing with repeated social interactions of this kind is to adopt a simple heuristic strategy in which one bases one's plays not on how one expects others to behave but rather on how they have behaved in the past. The best known of these heuristic strategies is TIT FOR TAT, which is composed of the following two rules:

- 1 Always cooperate in the first encounter
- 2 In any subsequent encounter do what your opponent did in the previous round

Theorists have found TIT FOR TAT a potentially powerful explanatory tool in explaining the evolutionary emergence of altruistic behavior for two reasons. First, it is simple and involved no complicated calculations. And second, it is what evolutionary game theorists call an *evolutionarily stable strategy* – that is to say, a population where there are sufficiently many “players” following the TIT FOR TAT strategy with a sufficiently high probability of encountering each other regularly will not be invaded by a subpopulation playing another strategy (such as the strategy of always defecting). TIT FOR TAT, therefore, combines simplicity with robustness.

Here, finally, we get to the cheater detection module. Simple though TIT FOR TAT is, it is not totally trivial to apply. It requires being able to identify instances of cooperation and defection. It involves being able to tell when an agent has taken a benefit without paying the corresponding price. An agent who consistently misidentifies defectors and free riders as cooperators (or, for that matter, vice versa) will not flourish.

This is why, according to Cosmides and Tooby, we evolved a specialized module in order to allow us to navigate social situations that depend crucially upon the ability to identify defectors and free riders. Since the detection of cheaters and free riders is essentially a matter of identifying when a conditional obligation has been breached, this explains why we are so much better at deontic versions of the selection task than ordinary versions – and why we are better, more generally, at conditional reasoning about rules, obligations, and entitlements than we are at abstract conditional reasoning.



**Exercise 8.1** Explain the argument from the evolution of cooperation to the cheater detection module in your own words.



### BOX 8.1 The Prisoner's Dilemma

A prisoner's dilemma is a strategic interaction that has the following puzzling and undesirable feature. If each participant does what seems to be the rational thing from their individual perspective, then the result for everyone is much worse than they could have achieved by cooperating.

The problem derives its name from a scenario where two prisoners are being separately interrogated by a police chief who is convinced of their guilt, but lacks conclusive evidence. He proposes to each of them that they betray the other, and explains the possible consequences. If each prisoner betrays the other then they will both end up with a sentence of 5 years in prison. If neither betrays the other, then they will each be convicted of a lesser offense and both end up with a sentence of 2 years in prison. If either prisoner betrays the other without himself being betrayed, however, then he will go free while the other receives 10 years in prison. Here is the pay-off table.

		PLAYER B	
		Betray	Not betray
PLAYER A	Betray	À5, À5	0, À10
	Not betray	À10, 0	À2, À2

Each entry represents the outcome of a different combination of strategies on the part of Prisoners A and B. The outcomes are given in terms of the number of years in prison that will ensue for Prisoners A and B, respectively (presented as a negative number, since years in prison are undesirable). So, the outcome in the bottom left-hand box is 10 years in prison for Prisoner A and none for Prisoner B, which occurs when Prisoner A does not betray, but Prisoner B does.

Imagine looking at the pay-off table from Prisoner A's point of view. You might reason like this. Prisoner B can do one of two things – betray me or not. Suppose he betrays me. Then I have a choice between 5 years in prison if I also betray him – or 10 years if I keep quiet. So, my best strategy if he betrays me is to betray him. But what if he does not betray? Then I have got a choice between 2 years if I keep quiet as well – or going free if I betray him. So, my best strategy if he is silent is to betray him. Whatever he does, therefore, I'm better off betraying him.

A game theorist would say that betray is Prisoner A's *dominant strategy*. A dominant strategy is one that promises greater advantage to that individual than the other available strategies, irrespective of what the other player does.

Unfortunately, Prisoner B is no less rational than you are, and things look exactly the same from her point of view. Her dominant strategy is also betray. So, you and Prisoner B will end up betraying each other and spending 5 years each in prison, even though you both would have been better off keeping silent and spending 2 years each in prison.



## Two Arguments

The extended case study of the cheater detection module is reinforced by two more general arguments, to which we turn now.

The arguments rest on two assumptions about evolution. The basic assumptions (surely both correct) are that the human mind is the product of evolution, and that evolution works by natural selection. These two basic assumptions give us a fundamental constraint upon possible mental architectures. Any mental architecture that we have today must have evolved because it was able to solve the adaptive problems that our ancestors encountered. Conversely, if you can show that a particular mental architecture could not have solved those adaptive problems, then it could not possibly be the architecture that we now have – it would have died out long ago in the course of natural selection.

In this spirit, the two arguments set out to show that evolution could not have selected a domain-general mental architecture. No domain-general, central processing system of the type that Fodor envisages could have been selected, because no such processing system could have solved the type of adaptive problems that fixed the evolution of the human mind.

**The argument from error.** This argument starts from the basic fact that what natural selection selects for are heritable traits that preserve fitness. But what counts as fitness? What are the criteria for fitness?

According to Cosmides and Tooby, these fitness criteria have to be domain-specific, not domain-general. What counts as fitness-promoting behavior varies from domain to domain. They give the example of how one treats one's family members. It is certainly not fitness-promoting to have sex with close family members. But, in contrast, it is fitness-promoting to help family members in many other circumstances. But not in every circumstance. If one is in a social exchange with a prisoner's dilemma-type structure and is applying something like the TIT FOR TAT algorithm, then it is only fitness-promoting to help family members who are cooperating – not the ones that are taking the benefit without paying the costs.

So, because there are no domain-general fitness criteria, there cannot (they argue) be domain-general cognitive mechanisms. Domain-general cognitive mechanisms could not have been selected by natural selection because they would have made too many mistakes – whatever criteria of success and failure they had built into them would have worked in some cases, but failed in many more. Instead, say Cosmides and Tooby, there must be a distinct cognitive mechanism for every domain that has a different definition of what counts as a successful outcome.



### Exercise 8.2 State the argument from error in your own words and evaluate it.

**The argument from statistics and learning.** Like the previous argument, this argument focuses on problems in how domain-general cognitive systems can discover what fitness consists in. The principal problem is that the world has what Cosmides and Tooby describe as a “statistically recurrent domain-specific structure.” Certain features hold



**Figure 8.6** The evolutionary biologist W. D. Hamilton (1936–2000). Jeffrey Joy

with great regularity in some domains, but not in others. These are not the sort of things that a general-purpose cognitive mechanism could be expected to learn.

Their example is the model of kin selection proposed by the evolutionary biologist W. D. Hamilton (Figure 8.6). The problem of kin selection is the problem of explaining why certain organisms often pursue strategies that promote the reproductive success of their relatives, at the cost of their own reproductive success. This type of self-sacrificing behavior seems, on the face of it, to fly in the face of the theory of natural selection, since the self-sacrificing strategy seems to diminish the organism's fitness.

Hamilton's basic idea is that there are certain circumstances in which it can make good fitness-promoting sense for an individual to sacrifice herself for another individual. From an evolutionary point of view, fitness-promoting actions are ones that promote the spread of the agent's genes. And, Hamilton argued, there are circumstances where an act of self-sacrifice will help the individual's own genes to spread and thereby spread the kin selection gene. In particular, two conditions need to hold:

*Condition 1* The self-sacrificer must share a reasonable proportion of genes with the individual benefiting from the sacrifice.

*Condition 2* The individual benefiting from the sacrifice must share the gene that promotes kin selection.

What counts as a reasonable proportion? This is where Hamilton's famous kin selection equation comes in. According to Hamilton, kin selection genes will increase when the following inequality holds:

$$R_{xy} B_y > C_x$$

Here the  $x$  subscript refers to the self-sacrificer and the  $y$  subscript to the beneficiary of the sacrifice. The term  $R_{xy}$  is a measure of how related  $x$  and  $y$  are. The term  $C_x$  measures the reproductive cost of kin selection to  $x$ , while  $B_y$  measures the reproductive benefit to  $y$ . In English, therefore, Hamilton's kin selection equation says that kin selection genes will spread when the reproductive benefit to the recipient of the sacrifice, discounted by the recipient's degree of relatedness to the self-sacrificer, exceeds the reproductive cost to the self-sacrificer.

Typically, two sisters will share 50 percent of their genes – or, more precisely, 50 percent of the variance in their genes (i.e., what remains after taking away all the genetic material likely to be shared by any two randomly chosen conspecifics). So, if  $x$  and  $y$  are sisters (and we measure relatedness in this way – evolutionary biologists sometimes use different measures), then we can take  $R_{xy} = 0.5$ . This tells us that it is only fitness-promoting for one sister to sacrifice her reproductive possibilities to help her sister when her sister will thereby do twice as well (reproductively speaking!) as she herself would have done if she hadn't sacrificed herself. So, the sacrifice will be fitness-promoting if, for example, the self-sacrificing sister could only have one more child, while the sacrifice enables her sister to have three more.

So much for the kin selection equation. Why should this make us believe in the massive modularity hypothesis? Cosmides and Tooby think that massive modularity is the only way of explaining how the kin selection law got embedded in the population. The kin selection equation exploits statistical relationships that completely outstrip the experience of any individual. According to Cosmides and Tooby, then, no domain-general learning mechanism could ever pick up on the statistical generalizations that underwrite Hamilton's kin selection law.

So how could the kin selection law get embedded in the population? The only way that this could occur, they think, is for natural selection to have selected a special-purpose kin selection module that has the kin selection law built into it.



### Exercise 8.3 State the argument from statistics and learning in your own words and evaluate it.

## Evaluating the Arguments for Massive Modularity

Both the argument from error and the argument from statistics and learning are compatible with the idea that human beings (not to mention other animals) are born with certain



innate bodies of domain-specific *knowledge*. This is a weaker requirement because information processing can exploit domain-specific knowledge without being modular.

Evolutionary psychologists are not always as precise as they could be in distinguishing between domain-specific modules and domain-specific bodies of knowledge. When we are thinking about the organization of the mind, however, the distinction is fundamentally important. When we formulate the massive modularity hypothesis in terms of cognitive modules it is a bold and provocative doctrine about the overall structure of the mind. It says that there is no such thing as a domain-general information-processing mechanism and that the mind is nothing over and above a collection of independent and quasi-autonomous cognitive subsystems.

But when we formulate the massive modularity thesis in terms of domain-specific bodies of knowledge it is much less clearly controversial. The idea that we (and quite possibly other animals) are born with innate bodies of knowledge dedicated to certain domains is not really a claim about the architecture of cognition. Cognitive scientists have proposed such innate bodies of knowledge in a number of different areas – such as numerical competence, intuitive mechanics, and so on.

But, on the other hand, even if one does not accept the massive modularity hypothesis in its strongest form, it still makes some very important points about the organization of the mind. In particular, it makes a case for thinking that the mind might be at least partially organized in terms of cognitive subsystems or modules that are domain-specific without having all the characteristics of full-fledged Fodorean modules. Cognitive scientists have taken this idea very seriously and we will be exploring it further in later chapters.

In Chapter 9 we will look at how the techniques of cognitive neuroscience can be used to study the organization of the mind, focusing in particular on the strengths and limits of using imaging techniques to map the mind. Chapters 13 and 14 develop a case study that brings the theoretical discussions about modularity to life. We will look at a debate that is very much at the forefront of contemporary cognitive science – the controversial question of whether there is a module responsible for reasoning about the mental states of others, or what many cognitive scientists have come to call the theory of mind module.

## 8.4

### Hybrid Architectures: The Example of ACT-R

This discussion of modules connects up with earlier discussions of information processing. Here's how. It may have occurred to you that the distinction between physical symbol systems and artificial neural networks is not all-or-nothing. Symbolic and distributed information processing seem to be suited for different tasks and for solving different types of problem.

The type of problems tackled by GOFAI physical symbol systems tend to be highly structured and sharply defined – playing checkers, for example, or constructing decision trees from databases. The type of problems for which artificial neural networks seem particularly well suited tend to be perceptual (distinguishing mines from rocks, for

example, or modeling how infants represent unseen objects) and involve recognizing patterns (such as patterns in forming the past tense of English verbs).

The extreme version of the physical symbol system hypothesis holds that *all* information processing involves manipulating and transforming physical symbol structures. It may be that Newell and Simon themselves had something like this in mind. There is a comparable version of the artificial neural networks approach, holding that physical symbol structures are completely redundant in modeling cognition – artificial neural networks are all we need. There seems to be room, though, for a more balanced approach that tries to incorporate both models of information processing. The ACT-R cognitive architecture developed by Michael Anderson and his research team at Carnegie Mellon University is a good example of how this might work.

The notion of a cognitive architecture, as used by computer scientists and psychologists, is a practical notion. A cognitive architecture is similar to a programming language. It gives researchers the tools to construct cognitive models using a common language and common tool kit.

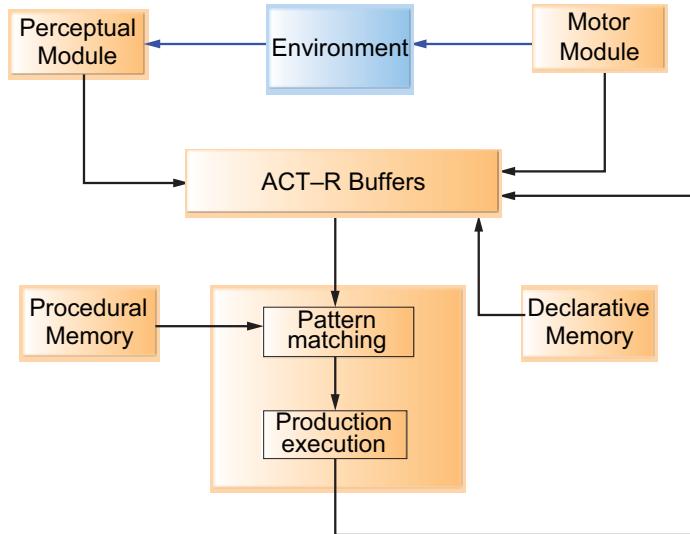
One of the first cognitive architectures was actually developed by Allen Newell, working with John Laird and Paul Rosenbloom. It was originally called SOAR (for “state operator and result”). The current incarnation is known as Soar. Soar is very closely tied to the physical symbol system hypothesis. It is based on the means–end and heuristic search approaches to problem solving that we looked at in Chapter 4. Soar is intended to be a unified model of cognition. It does not incorporate any elements corresponding to artificial neural networks. All knowledge is represented in the same way in the architecture, and manipulated in a rule-governed way.

But why not think about the mind in terms of *both* ways of modeling information processing, with some modules engaged in symbolic information processing and others in distributed information processing? That is the basic idea behind the ACT-R cognitive architecture developed by the psychologist John R. Anderson and his research team at Carnegie Mellon.

The ACT-R cognitive architecture is the latest installment of a cognitive architecture that was first announced under the name ACT in 1976 (“ACT” stands for “adaptive control of thought” and “R” for “rational”). It is a hybrid architecture because it incorporates both symbolic and subsymbolic information processing (partly in order to make it more neutrally plausible). One of the things that makes this architecture interesting from the perspective of this chapter is that it has a modular organization. Different modules performing different cognitive tasks and the type of information processing in a module depends upon the type of task the module performs.

## The ACT-R Architecture

The basic structure of ACT-R is illustrated in Figure 8.7. I want to draw your attention in particular to the two boxes at the top (the perceptual module and the motor module) and to the two modules directly below them (procedural memory and declarative memory).



**Figure 8.7** The ACT-R cognitive architecture.

In some versions of ACT-R the perceptual module is actually composed of a range of distinct modules – the visual module, the audition module, and so on, depending on the number of different types of input that the architecture can receive. Likewise, the motor module can sometimes itself be broken down into a speech module, a manual module, and so on, depending upon the different types of action that the architecture can perform. It is probably more helpful to talk generally about the *perceptual-motor layer*. Declarative and procedural memory collectively constitute what we can term the *cognitive layer*.

As we will see later, one very important feature of ACT-R is that communication between modules on different layers only takes place via *buffers*. A buffer is rather like a workspace. It contains the input that is available for processing by the relevant module. So, each module has its own buffer. Each perceptual module can only access sensory information that is in the relevant buffer (visual information in the visual buffer, and so on). Likewise, although not depicted in the figure, declarative and procedural memory each have their own buffer.

Another feature of ACT-R to note is that the two cognitive modules and the buffers from the modules in the perceptual–motor layer all feed into the pattern-matching and production execution mechanisms. This is where the actual decision-making takes place. We'll look at how it works further below.

For the moment I want to emphasize that the cognitive layer incorporates two fundamentally different types of knowledge – declarative and procedural. In philosophy this is often labeled the distinction between *knowledge-that* (declarative) and *knowledge-how* (procedural) – between, for example, knowing that Paris is the capital of France and knowing how to speak French.



**Exercise 8.4** Explain the distinction between knowledge-that and knowledge-how in your own words.

Declarative and procedural knowledge are both represented symbolically, but in different ways. Declarative knowledge is organized in terms of “chunks.” A chunk is an organized set of elements. These elements may be derived from the perceptual systems, or they may be further chunks. We can think of chunks as symbol structures (say the equation  $7 + 6 = 13$ ) built up in rule-governed ways from physical symbols (corresponding to **7**, **6**, **+**, **1**, and **3**). These chunks are stored in the declarative memory module.

ACT-R uses *production rules* to encode procedural knowledge. Production rules identify specific actions for the system to perform, depending upon which condition it finds itself in. When a production rule fires (as the jargon has it) in a given condition, it can perform one of a range of actions. It can retrieve a chunk from declarative memory, for example. Or it can modify that chunk – updating its representation of the environment, for example, or modifying a goal. It can also modify its environment. In this case the action really is an action – it sends a command to the motor module. And of course, production rules can be nested within each other, so that the output of a given production rule serves as a condition triggering the firing of another production rule. This allows complex abilities (such as multiplication) to be modeled as sets of production rules.

So far there is nothing hybrid about ACT-R. The way declarative and procedural knowledge is encoded and manipulated in the architecture is entirely in line with the physical symbol system hypothesis. And in fact, the same holds for the perceptual and motor modules. Here too information is encoded in the form of physical symbols. In some versions of ACT-R, the perceptual and motor modules are designed on the basis of the EPIC (Executive Process/ Interactive Control) architecture developed by David Kieras and David Meyer. EPIC falls squarely within the physical symbol system approach.

## ACT-R as a Hybrid Architecture

What makes ACT-R a hybrid architecture is that this symbolic, modular architecture is run on a subsymbolic base. Take another look at Figure 8.7. In many ways the overall organization looks very Fodorean. There are various modules, and they are all encapsulated. They communicate only via the buffer systems. And yet there is something missing. There is no system responsible for what Fodor would call central processing. But nor, on the other hand, is ACT-R massively modular. It does not have dedicated, domain-specific modules.

So, a natural question to ask of ACT-R is: How does it decide what to do? This is really the question of how it decides which production rules to apply? ACT-R is designed to operate serially. At any given moment, only one production rule can be active. But most of the time there are many different production rules that *could* be active. Only one of them is selected. How?

The job of selecting which production rule is to be active at a given moment is performed by the pattern-matching module. This module controls which production rule gains access to the buffer. It does this by working out which production rule has the highest utility at the moment of selection.



The production rule with the highest utility is the rule whose activation will best benefit the cognitive system. What counts as benefit depends upon the system's goals – or rather, to the system's current goal. The utility of a particular production rule is determined by two things. The first is how likely the system is to achieve its current goal if the production rule is activated. The second is the cost of activating the production rule.

So, the pattern-matching module essentially carries out a form of cost–benefit analysis in order to determine which production rule should gain access to the buffer. This cost–benefit calculation is really an application of the Bayesian approach discussed in Chapter 7. What the pattern-matching module is doing is based on calculating expected utility (and then factoring in the costs of applying the production rule).

The entire process takes place without any overseeing central system. It is a type of “winner-take-all” system. All the work is done by the equations that continually update the cost and utility functions. Once the numbers are in, the outcome is determined.

The designers of ACT-R describe these calculations as *subsymbolic*. This is a very important concept that is also standardly used to describe how artificial neural networks operate. For contrast, note that each production rule is purely symbolic. Production rules are built up in rule-governed ways from basic constituent symbols exactly as the physical symbol system hypothesis requires.

The compositional structure of production rules determines how the production rule behaves once it is activated, but it does not play a part in determining whether or not the rule is activated. For that we need to turn to the numbers that represent the production rule's utility. These numbers are subsymbolic because they do not reflect the symbolic structure of the production rule.

ACT-R has other subsymbolic dimensions, as summarized in Table 8.1. For example, it uses subsymbolic equations to model how accessible information is in declarative memory. The basic units of declarative memory are chunks – as opposed to the production rules that are the basic units of procedural memory. Each chunk has associated with it a particular activation level. This activation level can be represented numerically. The higher the activation level, the easier it is to retrieve the chunk from storage.

The activation levels of chunks in declarative memory are determined by equations. These equations are rather similar to the equations governing the utilities of production rules. There are two basic components determining a chunk's overall activation level. The first component has to do with how useful the chunk has been in the past. As before, usefulness is understood in terms of utility, which in turn is understood in terms of how the chunk has contributed to realizing the system's goals. The second component has to do with how relevant the chunk is to the current situation and context.

The example of ACT-R reveals two important lessons.

First, debates about the organization of the mind are closely connected to debates about the nature of information processing. Thinking properly about the modular organization of the mind requires thinking about how the different modules might execute their information-processing tasks. And second, different parts of a mental architecture might

**TABLE 8.1** Comparing the symbolic and subsymbolic dimensions of knowledge representation in the hybrid ACT-R architecture

	PERFORMANCE MECHANISMS		LEARNING MECHANISMS	
	SYMBOLIC	SUBSYMBOLIC	SYMBOLIC	SUBSYMBOLIC
Declarative chunks	Knowledge (usually facts) that can be directly verbalized	Relative activation of declarative chunks affects retrieval	Adding new declarative chunks to the set	Changing activation of declarative chunks and changing strength of links between chunks
Production rules	Knowledge for taking particular actions in particular situations	Relative utility of production rules affects choice	Adding new production rules to the set	Changing utility of production rules

exploit different models of information processing. Some tasks lend themselves to a symbolic approach. Others to a subsymbolic approach. The debate between models of information processing is not all-or-nothing.



## Summary

This chapter has focused on different ways of thinking about how the mind is organized. We began by looking at three different architectures for intelligent agents in AI, in order to see what distinguishes cognitive agents from simple reflex agents. Cognitive agents are standardly modeled in terms of quasi-autonomous information-processing systems, which raises the question of how those systems should be understood. Pursuing this question took us to Jerry Fodor's analysis of modular information-processing systems and explored his reasons for thinking that cognitive science is best suited to explaining modular systems, as opposed to nonmodular, central information-processing systems. We then examined an alternative proposed by massive modularity theorists, who hold that all information processing is modular. Finally, we turned to the hybrid architecture ACT-R, which brings the discussion of modularity into contact with the discussion of information processing in Part III. ACT-R is a modular system that combines the symbolic approach associated with the physical symbol system hypothesis and the subsymbolic neural networks approach.



## Checklist

**Computer scientists building intelligent agents distinguish different types of agent architectures.**

- (1) Simple reflex agents have condition-action rules (production rules) that directly link sensory and effector systems.
- (2) Simple reflex agents are not cognitive systems, unlike goal-based agents and learning agents.
- (3) Goal-based agents and learning agents are built up from subsystems that perform specific information-processing tasks.
- (4) This general approach to agent architecture raises theoretical questions explored in discussions of modularity.

### Fodor's Modularity Thesis

- (1) The thesis is built on a rejection of horizontal faculty psychology (the idea that the mind is organized in terms of faculties such as memory and attention that can process any type of information).
- (2) It proposes the existence of specialized information-processing modules that are: domain-specific informationally encapsulated mandatory fast
- (3) These modules may also have a fixed neural architecture and specific breakdown patterns.
- (4) Modules are employed for certain, basic types of information processing (e.g., shape analysis, color perception, and face recognition).
- (5) Modules provide inputs to nonmodular, central processing – the realm of belief fixation and practical decision-making, among other things.
- (6) Central processing is holistic and so not informationally encapsulated.

**According to the massive modularity hypothesis, all information processing is modular. There is no domain-general information processing.**

- (1) The human mind is claimed to be a collection of specialized modules, each of which evolved to solve a specific set of problems encountered by our Pleistocene ancestors.
- (2) Examples of these Darwinian modules are the cheater detection module and modules proposed for folk psychology (theory of mind) and folk physics (intuitive mechanics).
- (3) According to the argument from error, domain-general cognitive mechanisms could not have evolved because there are no domain-general fitness criteria.
- (4) According to the argument from statistics and learning, domain-general learning mechanisms cannot detect statistically recurrent domain-specific patterns (such as the kin selection equation proposed by W. D. Hamilton).
- (5) Both of these arguments can be satisfied with the much weaker claim that there are innate, domain-specific bodies of knowledge.

ACT-R is an example of a hybrid architecture that combines both symbolic and subsymbolic elements.

- (1) Knowledge in ACT-R is represented in two different ways – declarative knowledge is represented in chunks, while procedural knowledge is represented through production rules.
- (2) Items of knowledge become available for general information processing when they appear in one of the buffers. This general information processing is fundamentally symbolic in character.
- (3) In contrast, the processes that determine whether a particular item of knowledge ends up in a buffer are subsymbolic – equations, for example, that calculate how useful a given production rule might be in a particular context.
- (4) These processes are subsymbolic because they do not exploit or depend upon the internal symbolic structure of the item of knowledge.

## Further Reading

There is a useful introduction to intelligent agents in Russell and Norvig 2009, particularly chapter 2. An earlier version of this chapter (from the book's first edition) is available in the online resources. A good review can also be found in Poole and Mackworth 2010. See the online resources for other helpful collections pertaining to agent architectures.

Fodor's modularity thesis is presented in his short book *The Modularity of Mind* (1983). A summary of the book, together with peer commentaries, was published in the journal *Behavioral and Brain Sciences* (Fodor 1985). The summary is reprinted in Bermúdez 2006. For critical discussion of the modularity of face perception, see Kanwisher, McDermott, and Chun 1997 and Kanwisher 2000. Cosmides and Tooby have written an online evolutionary psychology primer, available in the online resources. More recent summaries of Cosmides and Tooby's research can be found in Cosmides, Barrett, and Tooby 2010 and Cosmides and Tooby 2013. Their 1994 paper discussed in the text is reprinted in Bermúdez 2006. It was originally published in Hirschfeld and Gelman 1994. This influential collection contains a number of other papers arguing for a modular approach to cognition. There is a useful entry on biological altruism in the online *Stanford Encyclopedia of Philosophy*. For Hamilton's theory of kin selection, see Dawkins 1979 (available in the online resources).

Pinker 1997 develops a view of the mind that integrates the massive modularity hypothesis with other areas of cognitive science. Pinker is a particular target of Fodor's discussion of massive modularity in Fodor 2000. Pinker responds to Fodor in Pinker 2005 (available in the online resources).

Carruthers 2006 is a book-length defense of a version of the massive modularity thesis. The journal *Mind and Language* published a precis of the book (Carruthers 2008b), together with three commentaries – Machery 2008, Wilson 2008, and Cowie 2008. Carruthers replies in the same issue (Carruthers 2008a). A good review of modularity research can be found in Barrett and Kurzban



2006. Also see Richard Samuels's chapter on massive modularity in Margolis, Samuels, and Stich 2012. The *Stanford Encyclopedia of Philosophy* also has an entry on modularity.

The homepage for the ACT architecture is the best place to start (see online resources). It contains a comprehensive bibliography with links to PDF versions of almost every referenced paper. For a brief overview of the general ACT approach, see Lebiere 2003. For a longer introduction to ACT-R, see Anderson et al. 2004. To see how ACT-R can be implemented neurally, see Zylberberg et al. 2011.





## CHAPTER NINE

# Strategies for Brain Mapping

<b>OVERVIEW</b>	229
<b>9.1 Structure and Function in the Brain</b>	230
Exploring Anatomical Connectivity	232
<b>9.2 Studying Cognitive Functioning: Techniques from Neuroscience</b>	237
Mapping the Brain's Electrical Activity: EEG and MEG	237
Mapping the Brain's Blood Flow and Blood Oxygen Levels: PET and fMRI	240
<b>9.3 Combining Resources I: The Locus of Selection Problem</b>	241
Combining ERPs and Single-Unit Recordings	242
<b>9.4 Combining Resources II: Networks for Attention</b>	246
Two Hypotheses about Visuospatial Attention	248
<b>9.5 From Data to Maps: Problems and Pitfalls</b>	249
From Blood Flow to Cognition?	250
Noise in the System?	251
Functional Connectivity versus Effective Connectivity	252



## Overview

This chapter explores what the wiring diagram of the mind looks like. This is a trickier question than it initially appears to be. Neuroanatomy (the study of the anatomical structure of the brain) is a good place to start, but neuroanatomy can only take us so far. We are looking for a *cognitive* wiring diagram. This takes us beyond anatomy, because cognitive functions rarely map cleanly onto brain areas. Section 9.1 looks in more detail at the theoretical and practical issues that arise when we start to think about the interplay between *structure* and *function* in the brain.

Many neuroscientists think that we can *localize* particular cognitive functions in specific brain areas (or networks of brain areas). Their confidence is in large part due to the existence of powerful techniques for studying patterns of cognitive activity in the brain. These techniques include

- EEG (electroencephalography) for measuring ERPs (event-related potentials)
- PET (positron emission tomography)
- fMRI (functional magnetic resonance imaging)



Section 9.2 introduces these techniques and their respective strengths, while the case studies in Sections 9.3 and 9.4 show how the different techniques can be combined to shed light on the complex phenomenon of attention.

Neuroimaging techniques do not provide a direct “window” on cognitive functions. They provide information about blood flow (in the case of PET) or the blood oxygen level dependent (BOLD) signal (in the case of fMRI). How we get from there to models of cognitive organization depends upon how we interpret the data. In Section 9.5 we will look at some of the challenges that this raises.



## 9.1

## Structure and Function in the Brain

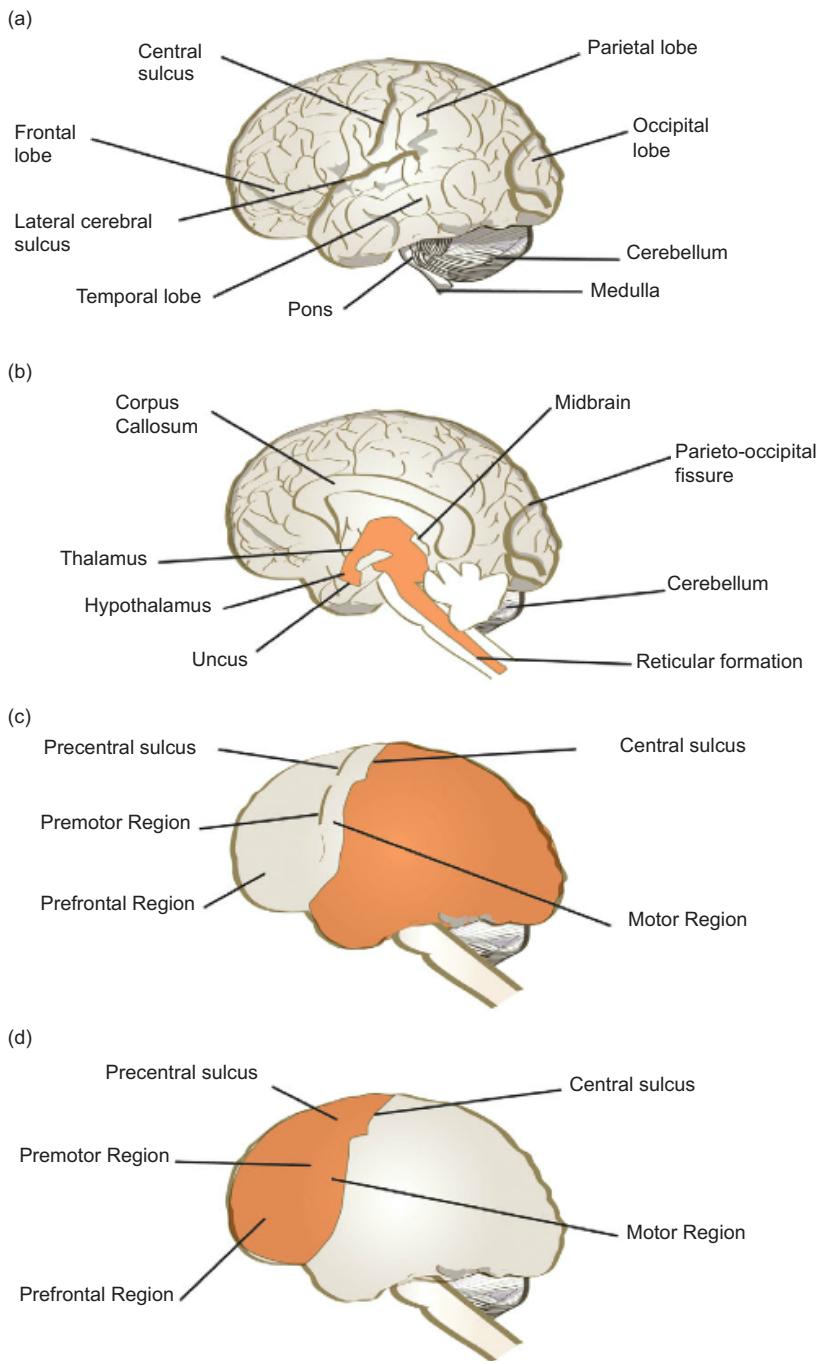
The brain has some conspicuous anatomical landmarks. Most obviously, it comes in two halves – the left hemisphere and the right hemisphere. The division between them goes lengthwise down the middle of the brain. As we saw in Section 3.2, the cortical surface of each of these hemispheres is divided into four lobes. Each of the four lobes is thought to be responsible for a different type of cognitive functioning. The frontal lobe is generally associated with reasoning, planning, and problem solving, for example. Anatomically speaking, however, the lobes are distinguished by large-scale topographic features known as *gyri* and *sulci* (singular: *gyrus* and *sulcus*).

If you look at a picture of the surface of the brain you will see many bumps and grooves. The bumps are the gyri and the grooves are the sulci. The sulci are also known as fissures. Many of these bumps and grooves have names. Some names are purely descriptive. The parieto-occipital sulcus, for example, separates the parietal lobe from the occipital lobe. Some of the names are more interesting. The Sylvian sulcus (which is marked in Figure 9.1 as the lateral cerebral sulcus) divides the temporal lobe from the lobe in front of it (the frontal lobe) and from the lobe above it (the parietal lobe). It is named after Franciscus Sylvius, who was a seventeenth-century professor of medicine at the University of Leiden in the Netherlands.

The diagram in Figure 9.1a is drawn from a review article published in *Scientific American* in 1970 by the famous Russian neuropsychologist Alexander Luria. It illustrates some of the most prominent large-scale features of the anatomy of the brain’s surface. My main interest in reproducing it, however, is to contrast it with the other three diagrams in Figure 9.1. Each of these depicts one of what Luria thought of as the three main *functional networks* in the brain. Luria called these networks “blocks.” They are colored brown in the diagrams.

According to Luria, each block has very different roles and responsibilities. Figure 9.1b is the most primitive block, made up of the brain stem and the oldest parts of the cortex. (This would be a good moment to look back at the first few paragraphs of Section 3.2.) According to Luria, this system regulates how awake and responsive we are. The second block (in Figure 9.1c) regulates how we code, control, and store information, while the third block (Figure 9.1d) is responsible for intentions and planning.

The specific details of Luria’s analysis are not particularly important. What he was reviewing in 1970 is no longer state of the art now. We are looking at Luria’s diagram



**Figure 9.1** Luria's (1970) diagram of the functional organization of the brain. The top diagram is anatomical, while the other three depict functional networks. (Adapted from Luria 1970)



because it clearly illustrates two things. The first is the difference between anatomy and cognitive function. The diagram in Figure 9.1a is an anatomical diagram. It organizes the brain in terms of large-scale anatomical features (such as the lobes and the sulci). It divides the brain into regions, but it has nothing to say about what those regions actually do.

The other three diagrams, however, are not purely anatomical. They mark many of the same anatomical regions, but they are organized in functional terms. This is particularly clear in Figures 9.1c and 9.1d, corresponding to Luria's second and third blocks. Here we have regions picked out in terms of what they are thought to do (in terms of the cognitive function that they serve). So, for example, a particular section of the frontal lobe is identified as the motor region (responsible for planning voluntary movements).

The second thing that we learn from Luria's diagram is how easy it is to slide from talking about anatomical areas to talking about functional areas (and vice versa). When we talk about the Sylvian sulcus we are talking about an anatomical feature of the brain. When we talk about the motor region, in contrast, we are talking about a region of the brain identified in terms of its function. But it is very common to have (as we have here) diagrams and maps of the brain that use both types of label. And in fact, the same area can have two very different names depending on how we are thinking about it. The precentral gyrus, for example, is an anatomical feature located just in front of the central sulcus. It is also called the primary motor cortex, because neuroscientists have discovered that directly stimulating this area causes various parts of the body to move.

## Exploring Anatomical Connectivity

One of the most fundamental principles of neuroscience is the *principle of segregation*. This is the idea that the cerebral cortex is divided into segregated areas with distinct neuronal populations. Neuroscientists still use a classification of anatomical areas in the cerebral cortex developed by the great German neuroanatomist Korbinian Brodmann in the late nineteenth and early twentieth century.

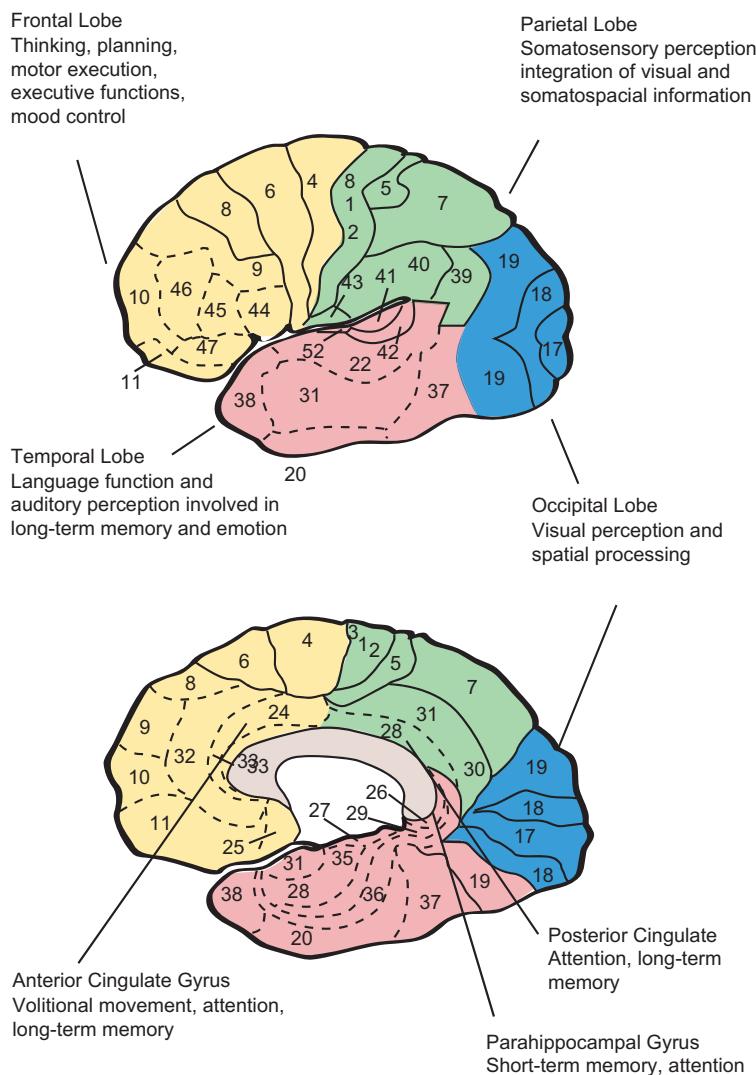
Brodmann's basic proposal was to distinguish different regions in the cerebral cortex in terms of the types of cell that they contain and how densely those cells occur. To study how cells are distributed in the cortex, Brodmann used recently discovered techniques for staining cells. Staining methods are still used by neuroscientists today. They involve dipping very thin slices of brain tissue into solutions that allow details of cellular structure to be seen under a microscope. Brodmann used the Nissl stain, developed by the German neuropathologist Franz Nissl. The Nissl stain turns all cell bodies a bright violet color.

By using the Nissl stain to examine the distribution of different types of neuron across the cerebral cortex, Brodmann identified over fifty different cortical regions. Figure 9.2 gives two views of the brain with the four lobes and the different Brodmann areas clearly marked. The top view is a lateral view (from the side) while the lower one is a medial view (down the middle).

Remarkably, Brodmann's classification of cortical regions also serves as a basis for classifying cortical regions according to their function. Here are some examples.



## Brodmann Areas



**Figure 9.2** Map of the anatomy of the brain showing the four lobes and the Brodmann areas. The captions indicate general functional specializations. The top view is a lateral view (from the side), while the lower one is a medial view (down the middle). (Reproduced courtesy of [www.appliedneuroscience.com](http://www.appliedneuroscience.com))



- The primary visual cortex, also known as area V1, is the point of arrival for information from the retina. In anatomical terms it is Brodmann area 17.
- Somatosensory information about the body gained through touch and body sense arrives in a region of the postcentral gyrus known as the primary somatosensory cortex. This is Brodmann area 3.
- We have already mentioned the primary motor cortex (the precentral gyrus). This is Brodmann area 4.

Even from an anatomical point of view, identifying segregated and distinct cortical regions can only be part of the story. We also need to know how the cortical regions are connected with each other. This would give an anatomical wiring diagram of the brain – or, to use the standard terminology, a map of *anatomical connectivity*.

Exploring anatomical connectivity requires a whole new set of techniques. One very influential technique is called *tract tracing*. This involves injecting a marker chemical into a particular brain region. Typical markers are radioactive amino acids or chemicals such as horseradish peroxidase (HRP). When the marker is injected near to the body of a nerve cell it is absorbed by the cell body and then transported along the cell's axon. Looking to see where the marker ends up allows neuroanatomists to identify where the cell projects to – and doing this for enough cells allows them to work out the connections between different brain regions.

Tract tracing is what is standardly called an invasive technique. It is only possible to discover where HRP has been transported to by examining sections of the cortex through a microscope. This cannot be done on living creatures. And so neuroanatomists have primarily worked on the brains of nonhuman animals – primarily macaque monkeys, rats, and cats. Their results are often represented using *connectivity matrices*.

Figure 9.3 is an example, from a very influential set of data on the visual system of the macaque monkey published in 1991 by Daniel J. Felleman and David Van Essen. The brain regions are abbreviated in a standard way. We can read off the matrix the regions to which any given region projects. Find the region you are interested in on the first column and then work your way across. If there is a “1” in the column corresponding to another brain region, then there is a connection going from the first to the second. If there is a “0” then no connection has been found. The gaps in the matrix indicate a lack of information.

The same data can be presented in a form that makes it look much more like a wiring diagram. We see this in Figure 9.4. The wiring diagram format makes it a little easier to visualize what is going on, but it doesn't give quite as much information as the connectivity matrix.



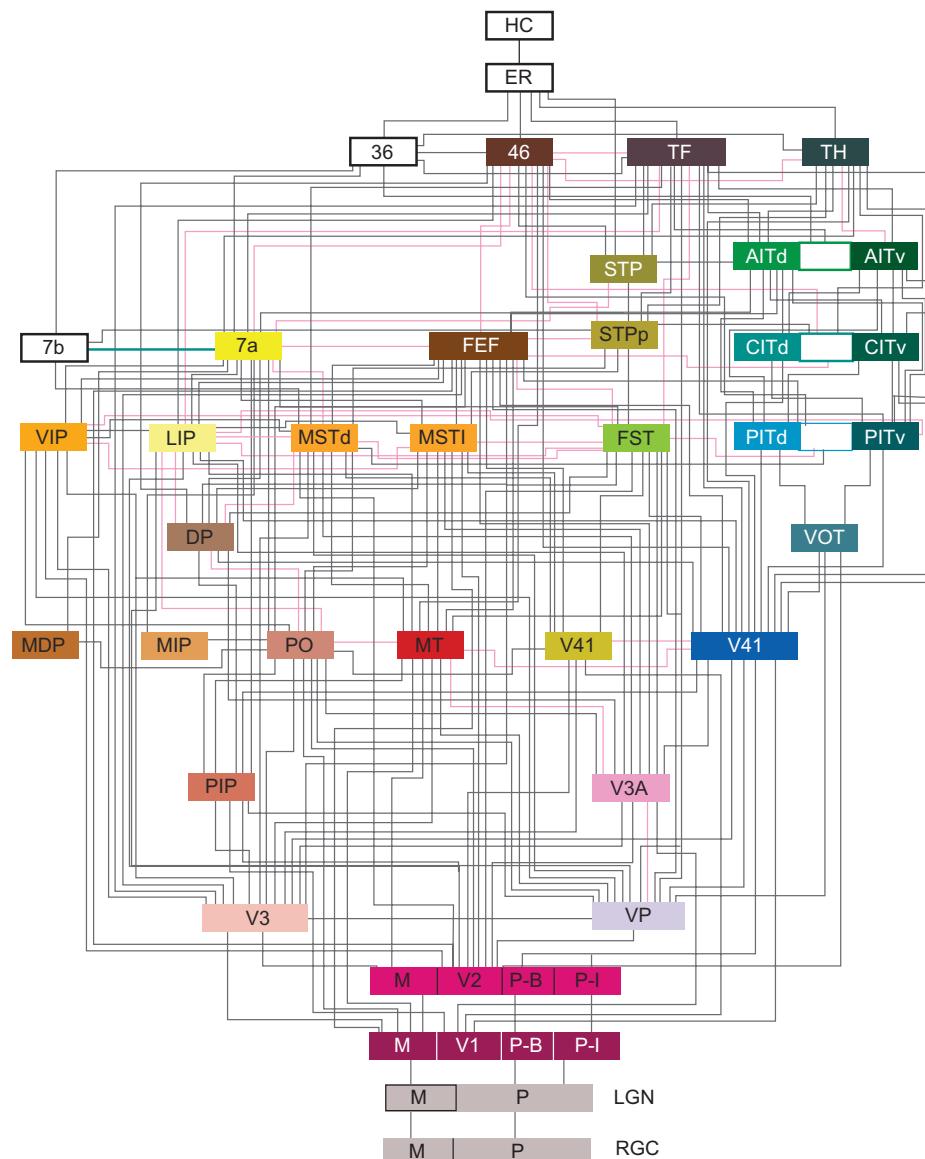
### Exercise 9.1 What type of information about anatomical connectivity do we get from a connectivity matrix but not from a wiring diagram?

Unfortunately, there are important limitations on what we can learn from information about anatomical connectivity.

- Anatomical connectivity data are largely derived from animal studies, whereas the brains that we are really interested in are our own. The information that we have about anatomical connectivity specifically in humans is largely derived from postmortem studies

	V1	V2	V3	VP	V3A	V4	VOT	V4T	MT	FST	PITd	PIT	PITv	CITd	CIT	CITv	AITd	AITv	STPp	STP	STPa	TF	TH	MStd	MSt1	PO	PP	LP	VIP	MIP	MDP	DP	7a	FEF	46
V1	1	1	0	1	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	0	0	0			
V2	1	1	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	0	0	0			
V3	1	1	1	1	1	1	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1	1	1	1	0	0	0			
VP	0	1	1	1	1	1	1	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1	1	1	1	0	0	0			
V3A	1	1	1	1	1	1	1	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	0	0	0	0			
V4	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	0	0	1	1	0	0	1	1	1	0	0	0	0			
VOT	0	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0	0	0			
V4t	1	1	0	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	0	0	0	0			
MT	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0		
FST	0	0	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	0		
PITd	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
PIT																																			
PITv	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
CITd	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
CIT																																			
CITv	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
AITd	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
AITv	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
STPp	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
STP																																			
STPa	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
TF	0	0	1	1	0	1	0	0	1	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
TH	0	0	0	0	0	0	1	0	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1		
MStd	0	1	1	1	1	1	0	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1		
MSt1	0	1	0	0	1	0	0	0	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0	0	0	1	0	0	1		
PO	1	1	1	1	1	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1		
PIP	1	1	1	1	1	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1		
LIP	0	1	1	1	1	1	1	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1	1	1	1	1	1	1	1		
VIP	0	1	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1		
MIP	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
MDP	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
DP	0	0	0	0	0	1	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1		
7a	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1	1		
FEF	0																																		
46	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1	1		

**Figure 9.3** A connectivity matrix for the visual system of the macaque monkey. (Adapted from Felleman and Van Essen 1991)



**Figure 9.4** An anatomical wiring diagram of the visual system of the macaque monkey. (Adapted from Fellman and Van Essen 1991)

of human brains. Techniques for studying human anatomical connectivity *in vivo*, such as *diffusion tractography*, are being developed, but are still in their infancy.

- Anatomical wiring diagrams do not carry any information about the direction of information flow between and across neural regions. There are typically at least as many feedback connections as feedforward connections.



- Anatomical connectivity is studied almost completely independently of cognitive functioning. An anatomical wiring diagram tells us which brain regions are in principle able to “talk” directly to each other. But it does not tell us anything about how different brain regions might form circuits or networks to perform particular information-processing tasks.



9.2

## Studying Cognitive Functioning: Techniques from Neuroscience

In addition to the principle of segregation, most neuroscientists also accept a *principle of integration*. This is the idea that cognitive functioning involves the coordinated activity of networks of different brain areas, with different types of task recruiting different networks of brain areas.

So, in order to make further progress in understanding how cognition works we need to supplement information about anatomical connectivity with information about what actually goes on in the brain when it is performing specific cognitive tasks. Neuroscientists have developed a number of techniques for doing this.

### Mapping the Brain's Electrical Activity: EEG and MEG

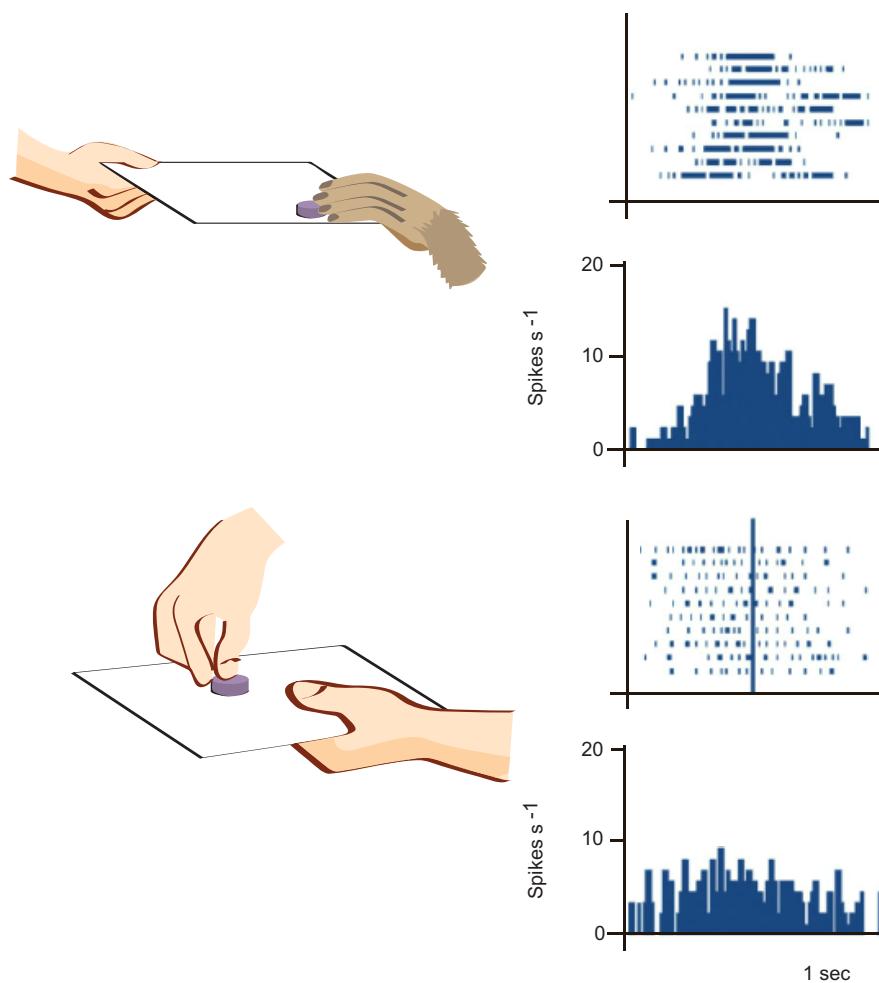
When neurons fire they send electrical impulses down their axons. These *action potentials* are transmitted to the dendrites of other neurons at *synapses*. Electrical synapses transmit electrical signals directly, while chemical synapses transmit chemicals called neurotransmitters. The precise details of how this works are not important for now. Two things are important. The first is that this electrical activity is a good index of activity in neurons. What neurons do is fire, and when they fire they generate electricity. The second is that there is a range of different techniques for measuring this activity.

Neurophysiologists can record the discharge of action potentials in individual neurons by placing a microelectrode close to the cell being recorded. (For an illustration see Figure 3.13.) This technique has been used to identify neurons that are sensitive to particular stimuli. The recent discovery of what are known as *mirror neurons* is a very good example. A group of neuroscientists led by Giacomo Rizzolatti in Parma, Italy, have identified neurons in monkeys that fire both when the monkey performs a specific action and when it observes that action being performed by an observer. This is illustrated in Figure 9.5.

To study the brain's organization and connectivity, however, we need to look at the electrical activity of populations of neurons, rather than single neurons.

Human encephalography (EEG) is a relatively straightforward tool for studying the activity of larger populations of neurons, requiring little complicated machinery or disturbance to the subject. It uses electrodes attached to the skull and wired up to a computer. Each electrode is sensitive to the electrical activity of thousands of neurons, with the neurons nearest the electrode making the largest contribution to the output signal.

The coordinated activity of these neural populations can be seen in EEGs as oscillatory waves at different frequencies. These frequencies are typically labeled in terms of bands.



**Figure 9.5** The results of single-neuron recordings of a mirror neuron in area F5 of the macaque inferior frontal cortex. The neuron fires both when the monkey grasps food (top) and when the monkey observes the experimenter grasping the food (bottom). Each horizontal line in the top diagram represents a single trial and each tick the firing of a neuron. Neural activity is summed over trials in the two histograms. (Adapted from Iacoboni and Dapretto 2006)

The bands are named with letters from the Greek alphabet – from alpha through to gamma. Confusingly, the alpha band is neither the lowest frequency nor the highest. The lowest frequency activity takes place in the delta band. Delta band activity is seen in very deep sleep (sometimes called slow wave sleep).



Name and example	Description
Delta	<p>Delta is the slow wave characteristic of deep, unconscious sleep. It is less than 4 Hz, and similar EEG frequencies appear in epileptic seizures and loss of consciousness, as well as some comatose states. It is therefore thought to reflect the brain of an unconscious person.</p> <p>The delta frequency tends to have the highest amplitude and the slowest frequency. Delta waves increase with decreasing awareness of the physical world.</p>
Theta	<p>Theta activity has a frequency of 3.5 to 7.5 Hz.</p> <p>Theta waves are thought to involve many neurons firing synchronously. Theta rhythms are observed during some sleep states, and in states of quiet focus, for example meditation. They are also manifested during some short-term memory tasks, and during memory retrieval.</p> <p>Theta waves seem to communicate between the hippocampus and neocortex in memory encoding and retrieval.</p>
Alpha	<p>Alpha waves range between 7.5 and 13 Hz and arise from synchronous (in-phase) electrical activity of large groups of neurons. They are also called Berger's waves in memory of the founder of EEG.</p> <p>Alpha waves are predominantly found in scalp recordings over the occipital lobe during periods of relaxation, with eyes closed but still awake. Conversely alpha waves are attenuated with open eyes as well as by drowsiness and sleep.</p>
Beta	<p>Beta activity is 'fast' irregular activity, at low voltage (12–25 Hz).</p> <p>Beta waves are associated with normal waking consciousness, often active, busy, or anxious thinking and active concentration.</p> <p>Beta is usually seen on both sides of the brain in symmetrical distribution and is most evident frontally. It may be absent or reduced in areas of cortical damage.</p>
Gamma	<p>Gamma generally ranges between 26 and 70 Hz, centered around 40 Hz.</p> <p>Gamma waves are thought to signal active exchange of information between cortical and other regions. They are seen during the conscious state and in REM dreams (Rapid Eye Movement Sleep). Note that gamma and beta activity may overlap in their typical frequency ranges, because there is still disagreement on the exact boundaries between these frequency bands.</p>

**Figure 9.6** Typical patterns of EEG waves, together with where/when they are typically found.  
(From Baars and Gage 2012)

In fact, different stages in the sleep cycle are associated with activity in different bands – and sleep specialists use EEG to identify and study sleep disorders. EEGs can be used for other forms of medical diagnosis. So, for example, epilepsy is associated with a distinctive, “spiky” wave, as can be seen in Figure 9.6.



EEGs are particularly important because they give a reliable way of measuring what are known as *event-related potentials* (ERPs). An ERP is the electrical activity provoked by a specific stimulus. ERPs are important because they have a very fine temporal resolution. In other words, they are sensitive to very small differences in elapsed time.

EEGs are not the only way of studying the electrical activity of large populations of neurons. But it is the most widespread technique and (not coincidentally, one imagines) the least expensive.

The other principal technology is *magnetoencephalography* (MEG). Magnetoencephalography measures the same electrical currents as are measured by EEG. It measures them through the magnetic fields that they produce. This allows a finer spatial resolution than is possible with EEGs. It is also much less susceptible to distortion due to the skull than EEG. But, on the other hand, it brings with it all sorts of technical issues. For example, it can only be carried out in a room specially constructed to block all alien magnetic influences, including the earth's magnetic field. MEG is relatively little used in research neuroscience (as opposed to medical diagnosis).

## Mapping the Brain's Blood Flow and Blood Oxygen Levels: PET and fMRI

We turn now to PET (*positron emission tomography*) and fMRI (*functional magnetic resonance imaging*), which were introduced in Sections 3.4 and 3.5, respectively. Instead of measuring electrical activity, these two techniques measure blood flow, since more blood flows to a particular brain region when it is active.

PET measures blood flow directly by tracking the movement of radioactive water. fMRI measures blood flow indirectly through blood oxygen levels in particular brain regions. This works because the increased neural activity in those regions does not consume all of the oxygen in the blood that reaches them. So, the ratio of oxyhemoglobin to deoxyhemoglobin increases in areas that see increased blood flow. This gives rise to the BOLD (*blood oxygen level dependent*) signal.

PET and fMRI are both much more sensitive to spatial change and variation than to change and variation over time. In this respect they are very different from EEG and MEG. Neuroimaging is much better at telling us about how cognitive activity is distributed across the brain over a period of time than they are at telling us about the precise sequence of events as information is processed.

This is acceptable, since functional neuroimaging is standardly used to identify networks of neural areas involved in carrying out cognitive tasks of a particular kind – those exploiting short-term memory, for example. This does not require a particularly fine temporal resolution. It simply requires being able to identify which neural regions are simultaneously active when the task is being performed. And the spatial resolution has to be sufficiently fine-grained for the results to be interpretable in terms of standard anatomical maps of the brain. The technology has to have sufficient spatial resolution to be able to pinpoint, for example, activity in the premotor cortex (Brodmann area 6), or in the orbitofrontal cortex (Brodmann area 11).

Table 9.1 summarizes some of the key features of these different techniques.

**TABLE 9.1** Comparing techniques for studying connectivity in the brain

	<b>DIRECTLY MEASURES</b>	<b>TEMPORAL RESOLUTION</b>	<b>SPATIAL RESOLUTION</b>
Single-unit recording	Potentials in individual neurons and very small populations of neurons	High	High
EEG (electroencephalography)	Electrical activity of larger populations of neurons	High	Low
MEG (magnetoencephalography)	Magnetic fields produced by electrical activity of larger populations of neurons	High	Low
PET (positron emission tomography)	Cerebral blood flow in particular brain regions	Low	High
fMRI (functional magnetic resonance imaging)	Levels of blood oxygen in particular brain regions	Low	High

As the table shows, there is no technique for studying larger populations of neurons that has both a high temporal resolution and a high spatial resolution. And so, as we'll see in the next two sections, neuroscientists have combined techniques in order to gain a more comprehensive perspective.

As the table shows, there is no technique for studying larger populations of neurons that has both a high temporal resolution and a high spatial resolution. Nonetheless, neuroscientists have combined techniques in order to gain a more comprehensive perspective. We will illustrate this through two case studies.

### 9.3

## Combining Resources I: The Locus of Selection Problem

Our first case study starts off from the basic fact that we experience the world in a highly selective way. At any given moment we effectively ignore a huge amount of the information that our perceptual systems give us. We saw an example of this in Chapter 1 – the so-called cocktail party phenomenon. The same holds for vision. In principle, we can see things that are more or less level with our ears. Yet we are barely aware of much of our so-called peripheral vision. It is only when something in the periphery “catches our eye” that we realize quite how far our field of vision extends.

This selectivity is a very basic feature of perception. We only *focus on* or *attend to* a small proportion of what we actually see, hear, touch, and so on. Psychologists label the mechanism responsible for this very general phenomenon *attention*.

In Chapter 1 we looked at Donald Broadbent's model of attention, in which attention functions as a selective filter. What the filter lets through depends upon what the cognitive system as a whole is trying to achieve. In a cocktail party situation, for example, the filter might be tuned to the sound of a particular individual's voice.



On Broadbent's model attention does its screening relatively early on in perceptual processing. The selective filter screens out all the sounds that don't correspond to the voice of the person I am talking to long before my auditory systems get to work on parsing the sounds into words and then working out what is being said. His model is what is known as an *early selection model*.

Other models claim that attention operates at a much later stage. These are *late selection models*. According to late selection models, important parts of perceptual processing are complete before attention comes into play. In vision, for example, late selection models hold that attention only comes into play once representations of sensory features (such as color, shape, and so on) have already been combined into representations of objects and those objects identified. The late selection approach is taken, for example, in the object-based model of attention developed by the cognitive psychologist John Duncan in the 1980s.



### Exercise 9.2 Explain in your own words the difference between late selection and early selection models of attention.

The *locus of selection problem* is the problem of determining whether attention is an early selection phenomenon or a late selection phenomenon. To solve it we need a way of tracking perceptual information processing to identify when attention comes into play.

## Combining ERPs and Single-Unit Recordings

The locus of selection problem is at bottom a problem about the temporal organization of information processing: Does the processing associated with selective attention take place before or after the processing associated with object recognition. This suggests using EEGs to measure the ERPs evoked by visual information processing. EEGs have a very high temporal resolution, sensitive at the level of milliseconds.

Remember that EEG (electroencephalography), is the general technique, while an ERP (the evoked reaction potential) is what the technique actually measures when it is *time-locked* with the onset of a particular stimulus. What we get from an ERP experiment is a wave that measures the electrical activity in the period of time immediately following the onset of the stimulus. The time is standardly measured in milliseconds (thousandths of a second), while the electrical activity is measured in microvolts (millionths of a volt).

We see a typical example in Figure 9.7. The graph displaying the ERP as a number of spikes and troughs. These are known as the *components* of the ERP and represent voltage deflections. The voltage deflections are calculated relative to a prestimulus baseline of electrical activity – which might, for example, be derived by measuring electrical activity at the tip of the nose.

**WARNING:** Please bear in mind a very confusing feature of ERP graphs. The *y-axis* represents negative activations above positive ones. This is very counterintuitive because



it means that, when the line goes up the electrical activity is actually going down! And vice versa.

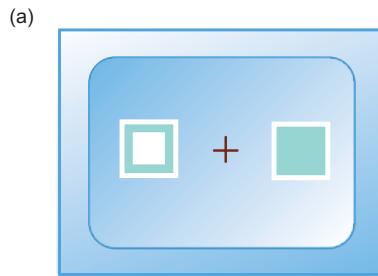
The time that elapses between stimulus onset and a particular spike or trough is known as the *latency* of the particular component. The components of the ERP for vision have been well studied. The earliest component is known as the C1 component. It is a negative component and appears at 50–90 ms after the appearance of the stimulus. There is a standard labeling for subsequent components. These are labeled either P or N, depending upon whether they are positive or negative. And they are given a number, which represents either their position in the ERP or their latency.

The P1 component, for example, is the first positive component, while the P300 is a positive component that occurs 300 ms (i.e., 0.3 seconds) after the stimulus is detected.

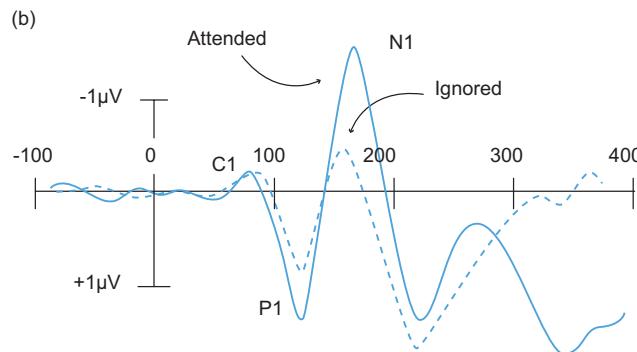
The P300 typically occurs in response to unexpected or novel stimuli. It is often taken as a sign that higher cognitive processes, such as attention, are involved in processing the stimulus. The graph in Figure 9.7b has the C1, N1, and P1 components marked. It is also possible to see the P200 component and a (slightly delayed) P300.

The ERP wave displays an *attention effect*. Certain components of the wave change depending upon whether or not the subject is attending to the stimulus. Figure 9.7a illustrates a typical experiment used to elicit the attention effect. The subject is asked to attend to one of two boxes in a screen. Stimuli are presented at various places in the screen and the ERPs are measured both for the case where the stimulus is in the box being attended to and the case where it is elsewhere.

The results of these experiments are illustrated in Figure 9.7b. The solid line shows the ERP when subjects are attending and the dotted line when subjects are not attending. There are important differences in two of the components – together with an important nondifference in one component. The nondifference first – there is no significant difference in the C1 component between the attended and the unattended cases. But there are significant differences in the P1 and N1 components. The P1 component is the first significant positive component and the N1 the first significant negative component. Both are larger when the subject is attending to the box in which the stimulus appears.



**Figure 9.7a** Common experimental design for neurophysiological studies of attention. The outline squares are continuously present and mark the two locations at which the solid square can be flashed. (Courtesy Stephen J. Luck and Michelle A. Ford)



**Figure 9.7b** Example of the occipital ERPs recorded in a paradigm of this nature. Note that the C1 wave (generated in area V1) shows no attention effect, whereas the P1 and the N1 waves (generated in extrastriate cortex) are larger for the attended stimuli. (Courtesy Stephen J. Luck and Michelle A. Ford)

This suggests that there are two additional bursts of information processing taking place roughly 100 and 200 ms after stimulus onset when attention is exercised. But how does it help us to decide whether attention is an early selection phenomenon or a late selection phenomenon?

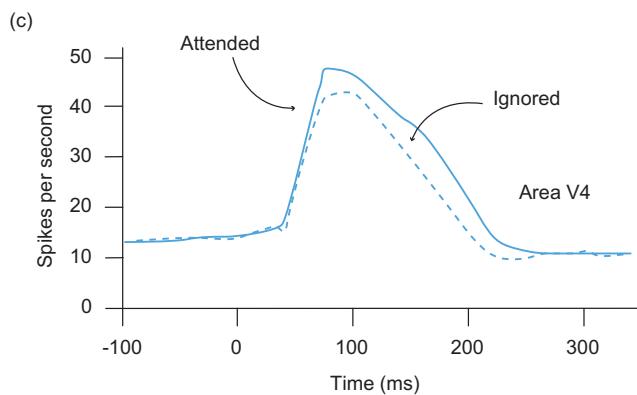
The ERP data on their own cannot settle this question. To make progress we need to triangulate ERP data with what we know about what goes on in different parts of the brain.

As we observed in Section 3.2, neurophysiologists generally accept that (at least in the visual system of the macaque monkey), object identification exploits the so-called ventral pathway that begins in V1 (the *striate cortex*) and then progresses through areas V2 and V4 en route to the inferotemporal cortex. Of these areas, V1 is responsible for processing basic shape. Visual areas V2 and V4 (which is an *extrastriate* area) are thought to process more advanced information about shape, together with information about color, texture, and so on.

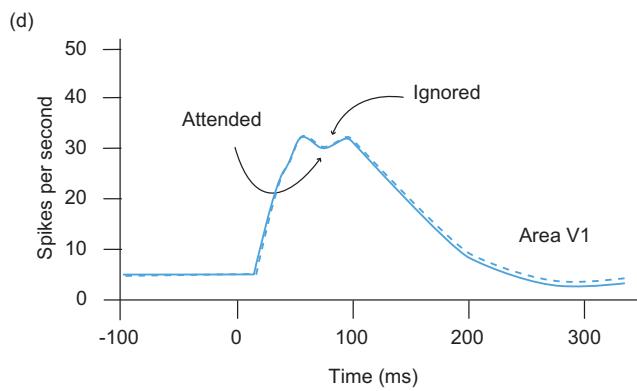
The different areas in the object identification pathway process different types of information separately but in parallel. Moreover, the information processing in V1, V2, and V4 is standardly thought to take place *upstream* of wherever the process of integrating all these different types of information takes place. In other words, all the information processing in the early visual areas such as V1, V2, and V4 takes place *before* the visual system is working with representations of objects.

This gives a clear criterion for thinking about the locus of selection problem. We said earlier that if attention is a late selection phenomenon then it only comes into play when the visual system has generated (and perhaps identified) representations of objects. Therefore, any evidence that the exercise of attention affects processing in the early visual areas will be evidence that attention is an early selection phenomenon.

This is why the ERP data are so significant. There is a range of evidence connecting different components of the ERP wave to processing in different visual areas. The C1



**Figure 9.7c** Single-unit responses from area V4 in a similar paradigm. Note that the response is larger for attended compared with ignored stimuli. (Courtesy Stephen J. Luck and Michelle A. Ford)



**Figure 9.7d** Single-unit responses from area V1 showing no effect of attention. (Courtesy Stephen J. Luck and Michelle A. Ford)

component, for example, is thought to reflect processing in the striate cortex (V1). Since the C1 component is constant across both the attended and the unattended conditions, we can conclude that processing in V1 is not modulated by attention. On the other hand, however, there is evidence connecting the P1 and N1 components with processing in the extrastriate cortex (i.e., in areas such as V2 and V4). So, the ERP data do seem to show that attention affects early visual processing, which supports the early selection rather than the late selection view.

This is reinforced by single-unit recordings. The diagrams in Figures 9.7c and 9.7d show the results of making recordings in areas V1 and V4 while monkeys are performing a task similar to that depicted in Figure 9.7a. As the graphs show, there is no difference between



levels of activity in V1 across the attended and unattended conditions. But there are significant differences in V4. This is certainly consistent with the hypothesis that attention is an early selection phenomenon.

There is a clear “take-home message” here. Although there are no techniques or technologies for studying cognitive activity directly and although each of the techniques has significant limitations, we can overcome many of the limitations by combining and triangulating the different techniques. The high temporal resolution of EEG complements the high spatial resolution of imaging technologies such as PET. And predictions from studies of humans using these techniques can be calibrated with electrophysiological studies on monkeys.



## 9.4

### Combining Resources II: Networks for Attention

Attention raises many important questions besides the locus of selection question. For example:

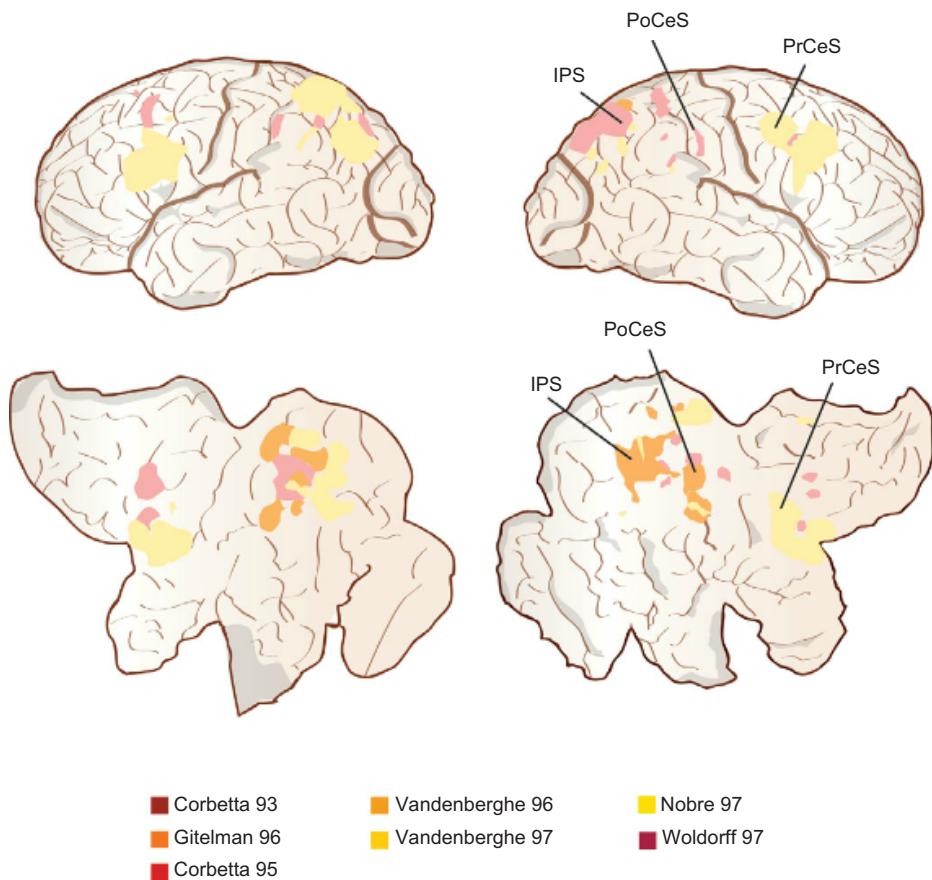
- Which brain areas are involved in attention?
- How is attention related to other cognitive processes, such as memory and action-planning?
- How does the brain direct attention to particular objects and particular places?

We will be exploring these questions in this section. This will allow us to see some of the power of experiments using functional neuroimaging – and also, to continue one of the themes of this chapter, to explore how neuroimaging data can be calibrated and reinforced with the results of electrophysiological experiments.

There are different varieties of attention. We can attend to one object among others – to the unfamiliar bird in the flock of sparrows, for example. Or we can attend to one part of an object rather than another – to the bird’s head or beak rather than its wings. Alternatively, we can attend to places – to the place where we expect the bird to fly to next.

The experiments that we looked at in the previous section focused on the last of these types of visual attention. Subjects were asked to focus on a particular location on the screen (marked by a box) – a location at which a stimulus might or might not appear. And it is what we will be focusing on in this section also. Neuroscientists and psychologists term it *spatially selective attention* (or *visuospatial attention*).

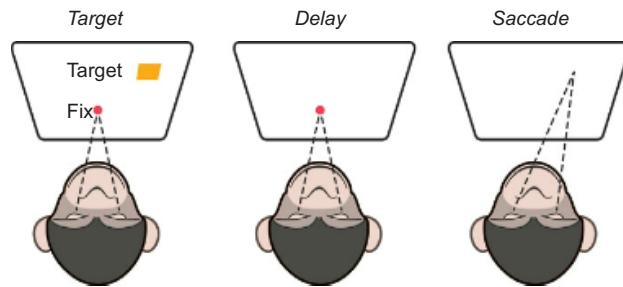
So, which brain areas are involved in spatially selective attention? Studies in the 1990s identified a network of cortical areas implicated in visuospatial attention. The specific tasks varied, but all of the experiments involved subjects directing attention to stimuli in the periphery of their visual field without moving their eyes. This is very important. Typically, we attend to different objects in the visual field by making very quick (and unconscious) eye movements known as *saccadic eye movements*. Experimenters studying visuospatial attention, however, are interested in attention as a mechanism that operates independently of eye movements – a mechanism that can be directed at different peripheral areas while gaze is fixated on a central point. Researchers call this *covert attention*.



**Figure 9.8** Frontoparietal cortical network during peripheral visual attention. Common regions of activation across studies include the intraparietal (IPS), postcentral (PoCeS), and precentral sulcus (PrCeS). (Adapted from Gazzaniga 2000)

Figure 9.8 summarizes a number of these studies. It identifies a network of areas in the parietal and frontal areas that are active during tasks that require subjects to direct covert attention to peripheral areas in the visual field. The existence of this frontoparietal cortical network is widely accepted among researchers into attention and has been confirmed by retrospective analyses of PET and fMRI data.

But simply identifying a network of brain areas involved in visuospatial attention does not in itself tell us much about how attention works. It does not tell us about how attention is related to other cognitive processes, such as memory or action-planning. And it does not tell us anything about how exactly the brain directs attention to particular locations in space. We turn to those questions now.



**Figure 9.9** An illustration of a typical delayed saccade task. The monkeys are trained to withhold their saccade to the visual target until the fixation point disappears. Note that the head does not move during the task. (From White and Snyder 2007)

## Two Hypotheses about Visuospatial Attention

There are two dominant hypotheses about how visuospatial attention works.

The first hypothesis is that visuospatial attention exploits certain memory mechanisms. The basic idea here is that, in order to attend to a specific location, we need actively to remember that location. If this is right, then we would expect brain networks associated with spatial working memory to be active during tasks that involve attention.

The second hypothesis is that there are very close connections between directing attention to a particular location and preparing to move to that location – even in the case of covert attention, where the intention to move is the intention to move the eyes. The prediction generated by this hypothesis is that brain areas associated with motor planning will be active in tasks that exploit visuospatial attention.

The two hypotheses are not necessarily exclusive. This is fortunate, because there is considerable experimental support for both of them.

Some of the evidence comes from single-neuron studies on monkeys. Carol Colby and her collaborators made recordings from an area in the parietal cortex known as LIP while monkeys were carrying out a *delayed saccade task*. LIP (the *lateral intraparietal* area – which we looked at in a different context in Section 7.3) is widely thought to play an important role in short-term memory of spatial locations.

In an ordinary saccade task the monkeys are trained to make a saccade (i.e., quickly move both eyes) from a central fixation point to a stimulus as soon as the stimulus appears. In a delayed saccade task the monkeys are trained not to make the saccade until the fixation point disappears – by which time the stimulus has disappeared (see Figure 9.9). When the fixation point disappears they then have to make a saccade to the location where the stimulus originally appeared.

Success on the delayed saccade task requires the monkeys to remember where the stimulus appeared if they are to make a successful saccade. This type of short-term memory about spatial location is typically called *spatial working memory*.



It turns out that the firing rates of neurons in LIP go up both when monkeys are performing delayed saccade tasks (and so exercising spatial working memory) and when they are carrying out peripheral attention tasks such as those discussed in the previous section.

This electrophysiological evidence is backed up by a wide range of neuroimaging studies on humans. Both PET and fMRI studies have shown significant overlap between the brain areas activated in visuospatial attention tasks and those active during tasks that require subjects to store and manipulate in working memory information about spatial locations. The results of these studies are depicted in the two diagrams on the left-hand side in Figure 9.10.

These diagrams show that, while there seem to be separate cortical networks for visuospatial attention and spatial working memory, these networks overlap very significantly in the parietal cortex. This is highly consistent with the results from the electrophysiological experiments.

Turning now to the relation between visuospatial attention and preparatory motor responses, the two diagrams on the right-hand side of Figure 9.10 report cross-experiment analyses. The experiments reported here all explored the relation between covert attention and saccadic eye movements. The diagrams superimpose the cortical networks thought to be involved in visuospatial attention onto the cortical networks implicated in saccadic eye movements.

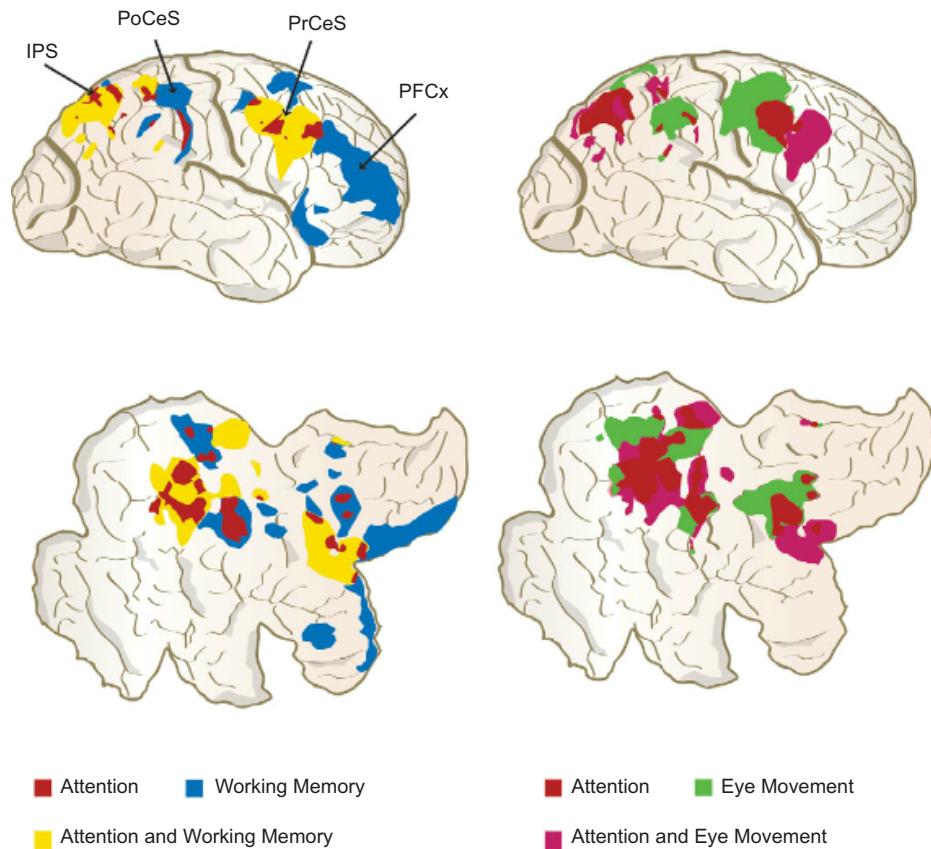
As an illustration, Maurizio Corbetta of Washington University in St. Louis scanned subjects both during conditions that required them to shift attention while maintaining their gaze fixed on a fixation point and during conditions in which gaze and attention shifted simultaneously. As the diagrams show, there is significant overlap across the covert attention and the saccadic eye movement tasks both in the parietal and in the precentral region (where the overlap is much stronger than in the working memory experiments).

This discussion of visuospatial attention illustrates an important methodological lesson. Progress in this area depends upon combining and calibrating what is learned from each of these techniques. We do not have any *direct* measures of cognitive activities such as visuospatial attention. But we do have the next best thing, which is a wide range of *indirect* measures. Single-unit recordings, PET, fMRI, and EEG all give us very different perspectives on visuospatial attention. We can use some techniques to compensate for the weaknesses of others. And we have powerful tools for cross-checking and integrating information from different sources. We have seen how this works in the case of visuospatial attention. This is an excellent case study in how neuroscientists are moving toward the goal of providing a cognitive wiring diagram of the brain.

## 9.5

## From Data to Maps: Problems and Pitfalls

Neuroimaging has yielded unparalleled insight into the structure and organization of the mind – perhaps more so than anything else in the neuroscientist's tool kit. But it is a tool that needs to be used with caution. Neuroimaging is *not* a direct picture of cognitive activity. It is easy to be seduced by the brightly colored images that emerge from software



**Figure 9.10** Peripheral attention versus spatial working memory versus saccadic eye movement across studies. Left: Regions active for peripheral attention (red), regions active for spatial working memory (blue), and regions of overlap (yellow). Note the striking overlap in parietal cortex, partial overlap in precentral region, and exclusive activation of prefrontal cortex (PFCx) for spatial working memory. Right: Comparison between peripheral attention (red) and saccadic eye movements (green). Note the strong overlap (magenta) in both parietal and precentral region. There is no activation in prefrontal cortex. (Adapted from Gazzaniga 2000)

packages for interpreting neuroimaging data. These images look very much like maps of the brain thinking. And so it is easy to think that neuroimaging gives us a “window on the mind.” In this section we will see why we need to be much more cautious.

### From Blood Flow to Cognition?

Neuroimaging only measures cognitive activity indirectly. fMRI measures the BOLD signal, while PET measures cerebral blood flow. But very little is known about how what we can



observe directly (the BOLD signal, for example) is connected to what we are trying to measure indirectly (information processing in the brain). As we saw in Section 3.6, there is a lively debate within neuroscience about the neural correlates of the BOLD signal. Researchers are calibrating fMRI data with electrophysiological techniques in order to try to work out whether the BOLD signal is correlated with the firing rates of populations of neurons, or whether it is correlated with the local field potentials (which are thought to reflect the inputs to neurons, rather than their outputs). We looked at some experimental evidence (from Logothetis and his collaborators) that seems to point to the second possibility.

But even if we had a conclusive answer to this question, we would still be a long way from a clear picture of the relation between variation in the BOLD signal and information processing in the brain. This is because we do not have any generally accepted models of how populations of neurons process information in particular brain areas. One illustration of this is that the BOLD signal gives us no indication whether the activity it measures is excitatory or inhibitory.

## Noise in the System?

Neuroimaging's great strength is its spatial resolution, but this comes at a cost.

The basic spatial unit in fMRI is the *voxel*. This is a three-dimensional version of a pixel (the name is a combination of the words “volume” and “pixel”). The basic unit of data obtained from fMRI is the BOLD signal in each voxel. The spatial resolution is directly correlated with the size of the voxels – the smaller the voxel, the higher the spatial resolution. The problem, though, is that the strength of the signal is directly correlated with the size of the voxel – the smaller the voxel, the lower the signal strength.

For some brain areas, particularly those involving basic perceptual processing or simple motor behaviors (such as finger tapping), experimenters can design tasks that elicit strong signals even when the voxel size is small. Things are not so straightforward, however, for more complex types of processing – particularly those performed by distributed networks of neural areas. Here it is often necessary to increase the voxel size in order to capture smaller fluctuations in the BOLD signal. Unsurprisingly, this decreases the spatial resolution. But it also has a less expected consequence.

Increasing the voxel size increases the range of different types of brain tissue occurring in each voxel. When the voxel includes extraneous material, such as white matter or cerebrospinal fluid, this can distort the signal, giving rise to what are known as *partial volume effects*. It can also happen that a single voxel contains more than one cell type, whereas neuroimaging data are standardly interpreted on the tacit assumption that voxels are homogeneous.

There are other ways in which noise can get into the system. Since everybody's brain is subtly different, meaningful comparisons across different subjects depend upon the data being *normalized* – that is, the data from each subject need to be reinterpreted on a brain atlas that uses a common coordinate system, or what is known as a *stereotactic map*. This



requires very complicated statistical techniques, which themselves may introduce distortion in the data.

And there are many different brain atlases, such as the Talairach-Tournoux atlas, the MNI atlas from the Montreal Institute of Neurology, and the Population-Average, Landmark and Surface-Based (PALS) atlas developed by David Van Essen at Washington University in St. Louis. Since different research groups often use a different atlas, this can make the business of comparing and contrasting different studies a tricky undertaking.

## Functional Connectivity versus Effective Connectivity

Neuroimaging helps us to understand the *connectivity* of the brain. But the wiring diagram that we get from fMRI and PET is still not quite the kind of diagram that we are looking for. Neither PET nor fMRI tells us anything *directly* about how information flows through the brain.

A single experiment can tell us which brain areas are simultaneously active while subjects are performing a particular task, but it does not tell us, for example, about the order in which the areas are active. The diagrams that present the results of neuroimaging experiments only show which areas “light up together.” They identify a network of areas that are simultaneously active when certain tasks are performed. But they really only identify correlations between the activity levels of different brain areas.

For this reason, neuroimagers distinguish *functional connectivity* from *effective connectivity*. Functional connectivity is a statistical notion. It is standardly defined in terms of statistical correlations between levels of activity in physically separate parts of the brain.

But to get a wiring diagram of how the brain works as an information-processing machine, we really need what neuroscientists call *effective connectivity*. Effective connectivity is a measure of how neural systems actually interact. Effective connectivity captures the idea that information processing is a causal process. Information flows through different brain areas in a particular order. What happens to the information at earlier stages affects how it is processed at later stages.



### **Exercise 9.3** Explain in your own words the distinction between anatomical, functional, and effective connectivity.

Neuroimaging is much better at telling us about functional connectivity than about effective connectivity. PET and fMRI are tools specialized for studying correlation, not causation.

Certainly, there are ways of deriving conclusions about effective connectivity from neuroimaging data. For example, one can design a series of experiments in a way that yields information about the flow of information. We looked at a very nice example of this back in Section 3.4. Steve Petersen and his collaborators were able to draw significant conclusions about the stages of lexical processing from a series of PET experiments using the paired-subtraction paradigm. The model that they developed is plainly a contribution to our understanding of the effective connectivity of the brain.



**Exercise 9.4** Look back at the lexical processing experiments described in Section 3.4 and explain in your own words how their experimental design overcomes some of the problems raised by the distinction between functional and effective connectivity.

And, as I have been stressing throughout this chapter, neuroimaging results can always be calibrated and triangulated with other tools and techniques, such as EEG and electrophysiology. Our discussion of the locus of selection problem showed how data from neuroimaging, EEG, and electrophysiology can be combined to develop a model of the effective connectivity of covert attention.

Nonetheless, we do have to be careful in how we interpret the results of neuroimaging experiments. In particular, we need to be very careful not to interpret experiments as telling us about effective connectivity when they are really only telling us about functional connectivity. We must be very careful not to draw conclusions about the causal relations between brain areas and how information flows between them from data that only tell us about correlations between BOLD signal levels in those areas.



## Summary

This chapter has looked at how cognitive neuroscience can help us to construct a wiring diagram for the mind. We began by highlighting the complex relations between functional structure and anatomical structure in the brain and then looked at some of the techniques for tracing anatomical connections between different brain areas. Completely different tools are required to move from anatomical connectivity to functional connectivity. We looked at various techniques for mapping the brain through measuring electrical activity and blood flow and blood oxygen levels. These techniques all operate at different degrees of temporal and spatial resolution. As we saw in two case studies, each having to do with a different aspect of the complex phenomenon of attention, mapping the functional structure of the brain requires combining and calibrating different techniques. At the end of the chapter we reviewed some of the pitfalls in interpreting neuroimaging data.

## Checklist

**It is a basic principle of neuroscience that the cerebral cortex is divided into segregated areas with distinct neuronal populations (the *principle of segregation*).**

- (1) These different regions are distinguished in terms of the types of cell they contain and the density of those cells. This can be studied using staining techniques.
- (2) This *anatomical* classification of neural areas can serve as a basis for classifying cortical regions according to their function.
- (3) Neuroscientists can study *anatomical connectivity* (i.e., develop an anatomical wiring diagram of the brain) by using techniques such as tract tracing or diffusion tractography.



- (4) Most of the evidence comes from animal studies. Neuroscientists have developed well worked out models of anatomical connectivity in macaque monkeys, rats, and cats.

**Neuroscientists also adopt the principle of integration – that cognitive functioning involves the coordinated activity of networks of different brain areas.**

- (1) Identifying these networks requires going beyond anatomical activity by studying what goes on in the brain when it is performing particular tasks.
- (2) Some of the techniques for studying the organization of the mind focus on the brain's electrical activity. These include electrophysiology, EEG, and MEG.
- (3) These techniques all have high temporal resolution – particularly EEG when it is used to measure ERPs. But the spatial resolution is lower (except for electrophysiology using microelectrodes).
- (4) Other techniques measure blood flow (PET) and levels of blood oxygen (fMRI). These techniques have high spatial resolution, but lower temporal resolution.

**The locus of selection problem is the problem of determining whether attention operates early in perceptual processing, or upon representations of objects. It provides a good illustration of how neuroscientists can combine different techniques.**

- (1) The problem has been studied using EEG to measure ERPs. Attentional effects appear relatively early in the ERP wave following the presentation of a visual stimulus.
- (2) These results can be calibrated with PET studies mapping stages in the ERP wave onto processing in particular brain areas. This calibration reveals attentional effects in areas such as V2 and V4, which carry out very basic processing of perceptual features.
- (3) This resolution of the locus of selection problem seems to be confirmed by single-unit recordings in monkeys.

**Neuroimaging techniques can help identify the neural circuits responsible for attention.**

- (1) Preliminary evidence from brain-damaged patients (e.g., with hemispatial neglect) points to the involvement of frontal and parietal areas in visuospatial attention.
- (2) This has been confirmed by many experiments on covert attention using PET and fMRI.
- (3) PET and fMRI experiments on humans, together with single-neuron experiments on monkeys, have shown that tasks involving visuospatial attention also generate activation in brain networks responsible for planning motor behavior and for spatial working memory.

**The discussion of attention shows that neuroimaging is a very powerful tool for studying cognition. It is not a “window on the mind,” however, and neuroimaging data should be interpreted with caution.**

- (1) Neuroimaging techniques can only measure cognitive activity indirectly. PET measures blood flow and fMRI measures the BOLD signal. There is a controversy in neuroscience about what type of neural activity is correlated with the BOLD signal (see Section 3.6) – and no worked out theory about how that neural activity functions to process information.
- (2) There are many opportunities for noise to get into the system in neuroimaging experiments. Partial volume effects can occur when the voxel size is large, and distortions can occur when data are being normalized to allow comparison across subjects.



- (3) Neuroimaging techniques are much better at telling us about *functional connectivity* (correlations between activation levels in different brain areas as a task is performed) than about *effective connectivity* (how information flows between different brain areas and how they influence each other).

## Further Reading

The explosion of interest in cognitive neuroscience in the last couple of decades has generated a huge literature. For keeping up to date with contemporary research, the journal *Trends in Cognitive Sciences* regularly contains accessible survey articles. Authoritative review articles on most of the key topics studied by cognitive neuroscientists can be found in the fifth edition of *The Cognitive Neurosciences*, edited by Michael Gazzaniga and George Mangun (2014). The four earlier editions (Gazzaniga 1995, 2000, 2004, 2009) also contain much useful material. Gazzaniga is one of the authors of an influential textbook on cognitive neuroscience (Gazzaniga, Ivry, and Mangun 2013). Chapter 3 is a useful introduction to the methods of cognitive neuroscience. Also see Baars and Gage 2010.

Zeki 1978 was one of the first papers to identify functional specialization in the primate visual system. David Van Essen's work is accessibly presented in Van Essen and Gallant 1994. The much-cited paper discussed in the text is Felleman and Van Essen 1991. Reviews of other classic work can be found in Colby and Goldberg 1999 and Melcher and Colby 2008. Orban, Van Essen, and Vanduffel 2004 is an interesting discussion of the challenges in comparing the neurobiology of cognitive function across humans and macaque monkeys. Also see Passingham 2009. An interesting trend in recent discussions of anatomical connectivity has been the use of mathematical tools from graph theory – in particular the idea of small-world networks. There is a very useful introduction in Bassett and Bullmore 2006. See Minati et al. 2013 for a more up-to-date review. Jirsa and McIntosh 2007 is a collection of surveys of different aspects of neural connectivity. For article-length surveys, see Ramnani et al. 2004, Bullmore and Sporns 2009, and Friston 2011. Bressler et al. 2008 uses Granger causality to explore effective connectivity in the neural basis of visual-spatial attention. For more on Granger causality, see Deshpande and Hu 2012 and Friston et al. 2013.

There has been much discussion of the pitfalls and advantages of using neuroimaging techniques to study cognitive function in the human mind. In addition to research on the neural basis of the BOLD signal discussed in Chapter 3 (see the references there), researchers have focused on the methodology of inferring cognitive function from selective patterns of activation. See, for example, Henson 2006 and Poldrack 2006. For a review of the state of fMRI at the time from a leading researcher, see Logothetis 2008. Also see Ashby 2011, Charpac and Stefanovic 2012, Machery 2012, and Poldrack, Mumford, and Nichols 2011. Poldrack 2018 is a book-length survey of the power and pitfalls of neuroimaging.

For surveys of research into selective attention, see Hopfinger, Luck, and Hillyard 2004 and Thigpen and Keil 2017. Experimental work reported in Section 9.3 is described more fully in Luck and Ford 1998. Stephen Luck is the author of an important textbook on ERP techniques (Luck 2005). The introductory chapter can be downloaded from the online resources. See also his coedited volume Luck and Kappenman 2011.

Humphreys, Duncan, and Treisman 1999 contains many useful papers on the psychology and neuroscience of attention, as does Posner 2004. For more details of the findings discussed in Section 9.4, see Chelazzi and Corbetta 2000. Other good reviews on a wide variety of attention phenomena can be found in chapters 8 and 10 of Baars and Gage 2010 as well as in Carrasco 2011, Chun, Golomb, and Turk-Browne 2011, Posner 2017, and Carrasco 2018.

**PART III**

**APPLICATIONS**









## CHAPTER TEN

# Models of Language Learning

### OVERVIEW 259

#### 10.1 Language and Rules 260

Understanding a Language and Learning a Language 261

#### 10.2 Language Learning and the Language of Thought: Fodor's Argument 263

#### 10.3 Language Learning in Neural Networks 266

#### The Challenge of Tense Learning 267

Connectionist Models of Tense Learning 269

#### 10.4 Bayesian Language Learning 274

Probabilities in Word and Phrase

Segmentation 275

Understanding Pronouns 276

Learning Linguistic Categories 278



## Overview

Language is a highly sophisticated cognitive achievement. Without it our cognitive, emotional, and social lives would be immeasurably impoverished. And it is a truly remarkable fact that almost all human children manage to arrive at more or less the same level of linguistic comprehension and language use. Unsurprisingly, cognitive scientists have devoted an enormous amount of research to trying to understand how languages are learned. This chapter looks at language learning from three of the theoretical perspectives discussed in earlier chapters:

- The language of thought hypothesis (a version of the physical symbol systems hypothesis)
- Connectionist neural networks
- Probabilistic Bayesian models

Section 10.1 introduces some of the basic theoretical challenges in explaining how we understand and learn languages. Since language is a paradigmatically rule-governed activity, it can seem very plausible to conceptualize linguistic understanding as a matter of deploying linguistic rules. This raises the question of where knowledge of the rules comes from. Answering that question is an important part of explaining how languages are learned.

We look at one answer to that question in Section 10.2. According to Jerry Fodor, young children learn linguistic rules by a process of hypothesis formation and testing. This process is itself a linguistic activity. According to Fodor, though, it cannot be carried out in a natural language. He

thinks that it takes place in the language of thought and he draws the natural conclusion that the language of thought must be innate.

Section 10.3 explores the very different connectionist approach, which uses neural networks to model how languages might be learned without explicitly representing rules. We look at how neural networks can be trained to learn the past tense of English verbs, both regular and irregular. As we'll see, the learning trajectory of these networks bears striking resemblances to the learning trajectory of human infants coming to terms with the complexities of verbs in English.

Like the connectionist approach, the Bayesian models that we consider in Section 10.4 have a positive account to offer of how languages are actually learned, often explicitly setting themselves against *nativist* or *innatist* accounts. Bayesians think that arguments for innatism about language have significantly underestimated how much can be learned through sensitivity to statistical regularities and through applying the type of Bayesian principles that we looked at in Chapter 7. We will look at three different examples of this general approach.

## 10.1

### Language and Rules

In many ways, speaking and understanding a natural language is the paradigm of a rule-governed activity. At a most basic level, every language is governed by grammatical rules. These rules, painfully familiar to anyone who has tried to learn a second language, govern how words can be put together to form meaningful sentences. But grammatical rules are only the tip of the iceberg. Linguists devote much of their time to trying to make explicit much more fundamental rules that govern how languages work. (These additional rules are more fundamental in the sense that they are supposed to apply to all languages, irrespective of the particular grammar of the language.)

Back in Section 1.3 we looked briefly at the version of transformational grammar proposed by Noam Chomsky in the 1950s. In effect, what Chomsky was proposing were rules that governed how a sentence with one type of grammatical structure could be legitimately transformed into a sentence with a different grammatical structure but a similar meaning.

The example we looked at there was the pair of sentences "John has hit the ball" and "The ball has been hit by John." Here we have two sentences with very different surface grammatical structures, but that convey similar messages in virtue of having the same deep (or phrase) structure. Chomsky's insight was that we can understand what is common to these sentences in terms of the rules that allow one to be transformed into the other. These are the transformational rules. Chomsky's view on what these rules actually are has changed many times over the years, but he has never abandoned the basic idea that the deep structure of language is governed by a body of basic rules.

The rule-governed nature of language makes thinking about language a very interesting test case for comparing and contrasting the different models of information processing that we looked at in Part II. These models take very different perspectives on the role of rules. As we saw in Chapter 4, the basic idea behind the physical symbol system hypothesis is that



information processing is a matter of manipulating physical symbol structures according to rules that are explicitly represented within the system. In contrast, in Chapter 5 we learned that it is not really possible to distinguish rules and representations in artificial neural networks (apart from the algorithm that governs how the network updates its activation levels). Information processing in artificial neural networks does not seem to involve rule-governed symbol manipulation. Nor, as we saw in Chapter 7, do Bayesian models explicitly encode rules (besides those required to manipulate probabilities and utilities).

Still, the fact that languages are governed by rules does not automatically mean that the information processing involved in understanding and learning languages has to involve manipulating symbol structures according to rules. If we are to arrive at that conclusion it will have to be through some combination of theoretical argument and empirical evidence. We turn now to some of the theoretical reasons that have been given for thinking that the physical symbol system hypothesis (particularly in its language of thought incarnation) is the only way of making sense of the complex phenomenon of linguistic comprehension and language learning. And then, in the rest of the chapter, we will test the power of those arguments by looking at neural network and Bayesian models of specific aspects of language learning.

## Understanding a Language and Learning a Language

What is it to understand a language? In a very general sense, there are two different dimensions to linguistic comprehension. One dimension is understanding what words mean. There is no language without vocabulary. But words on their own are not much use. The basic unit of communication is not the word, but rather the sentence. The logician and philosopher Gottlob Frege famously remarked that only in the context of a sentence do words have meaning. This takes us to the rules that govern how words can be put together to form meaningful sentences. As we have already seen, these rules are likely to fall into two groups. On the one hand there are the rules that tell us which combinations of words are grammatical. On the other there are the rules that govern the deep structure of language.

So, understanding a language is partly a matter of understanding what words mean, and partly a matter of understanding how words can be combined into sentences. What does this understanding consist in? The default hypothesis is that understanding a language is fundamentally a matter of mastering the relevant rules. This applies to the vocabulary of a language no less than to its grammar and deep structure. We can think of understanding the meaning of a word in terms of mastery of the rule that governs its application – the rule, for example, that the word “dog” refers to four-legged animals of the canine family and the rule that the word “square” applies to four-sided shapes with sides of equal size and each corner at an angle of 90 degrees.

The default hypothesis does not, however, tell us very much. Everything depends on how we think about mastering a rule. At one extreme is the view that there is no more to mastering a linguistic rule than being able to use words in accordance with the rule.

There is no need for competent language users to represent the rule in any way. All they need to be able to do is to distinguish applications of the word that fit the rule from applications that do not. This is a very minimalist conception of linguistic understanding, associated with some of the followers of the philosopher Ludwig Wittgenstein. It makes linguistic understanding much more of a practical ability than a theoretical achievement.

Many cognitive scientists, in contrast, think that this way of thinking about mastery of rules is far too weak. After all, the rock that falls in accordance with Newton's law of gravity cannot in any sense be said to have mastered that law. Mastering linguistic rules certainly requires using words in accordance with the rule, but it is not just a practical ability. Many theorists take the view that we cannot take linguistic abilities as given. They have to be explained in some way. And one explanation many have found plausible is that language users are capable of using words in accordance with linguistic rules because they represent those rules. These representations are thought to guide the language user's use of language. Language users use words in accordance with the rule because they somehow manage to compare possible sentences with their internalized representations of the rules. This is the other extreme. It makes linguistic understanding much more of a theoretical achievement than a practical ability – or rather, it takes linguistic understanding to be a practical ability grounded in a theoretical achievement.

So, the default hypothesis that linguistic understanding consists in mastery of linguistic rules can be understood in many different ways, depending on where one stands in relation to these two extremes. And this has significant implications for how one thinks about the information processing involved in understanding and using a language. The more importance one attaches to the explicit representation of rules, the more likely one is to think that this information processing must be understood through the physical symbol system hypothesis. This is because the physical symbol system hypothesis allows rules to be explicitly represented within the system.

Moreover, how one thinks about linguistic understanding has direct implications for how one thinks about language learning. The end point of language learning, after all, is linguistic understanding. This is why accounts of what it is to understand a language and what it is to learn a language tend to go hand in hand. We will be exploring this interdependence in the next few sections. In particular, we will be looking at issues of *innatism* or *nativism* (terms that I'll use interchangeably). Some accounts of linguistic understanding have the consequence that aspects of language must be innate, because they seem in principle to be incapable of being learned by young children, given the limited resources and evidence that children have.

We will look at one example of this type of theory in the next section, where we consider Fodor's argument that language learning requires an innate language of thought. And then in Sections 10.3 and 10.4 we will look at two prominent anti-innatist approaches – connectionist models of language in Section 10.3 and then Bayesian approaches in Section 10.4.



## 10.2

# Language Learning and the Language of Thought: Fodor's Argument

This section examines a powerful line of argument working backward from a strong conception of what it is to master linguistic rules to the conclusion that we should think about language learning in terms of the physical symbol system model – and, in particular, in terms of the language of thought hypothesis. This argument is due to the philosopher Jerry Fodor, although its basic thrust is, I think, one that many cognitive scientists would endorse and support.

Fodor starts off with a strong version of the rule-based conception of language learning. He thinks of the process of acquiring a language as a lengthy process of mastering the appropriate rules, starting with the simplest rules governing the meaning of everyday words, moving on to the simpler syntactic rules governing the formation of sentences, and then finally arriving at complex rules such as those allowing sentences to be embedded within further sentences and the complex transformational rules discussed by Chomsky and other theoretical linguists.

How does Fodor get from the rule-based conception of language learning to the existence of a language of thought? His argument is in his book *The Language of Thought*. It starts off from a particular way of thinking about the rules governing what words mean. According to Fodor these rules are what he calls *truth rules*. They are called truth rules because they spell out how words contribute to determining what it is for sentences in which they feature to be true. Mastering truth rules may not be all that there is to understanding a language. But Fodor is emphatic that we will not be able to understand a language without mastering truth rules. Truth rules may not be sufficient, but they are certainly necessary (he claims).

Let us take a very simple sentence to illustrate how truth rules work. Consider, for example, the sentence “Felicia is tall.” This sentence is what logicians call an atomic sentence. It is made up simply of a proper name (“Felicia”) and a predicate (“\_\_\_\_ is tall,” where the gap indicates that it needs to be “completed” by a name of some sort). Proper names are names of individuals and predicates are names of properties. And so, this gives us a very straightforward way of thinking about what makes an atomic sentence such as “Felicia is tall” true. The sentence is true just if the individual named by the proper name (i.e., Felicia) does indeed have the property named by the predicate (i.e., the property of being tall). So, the atomic sentence “Felicia is tall” is true just if Felicia is tall. It is standard to call this the *truth condition* of the sentence.

You may well think, though, that the truth condition cannot be much help to us in thinking about what it is to understand the sentence “Felicia is tall,” or about how one might learn how to use the expressions “Felicia” and “\_\_\_\_ is tall.” Here is the truth condition:

TC “Felicia is tall” is true just if Felicia is tall

Surely, you might say, someone can only understand the truth condition (TC) if they already understand the sentence “Felicia is tall” (because this very sentence features in the truth condition, both inside and outside quotation marks). But then the truth condition can only be intelligible to someone who already understands the expressions “Felicia” and “  is tall.” It cannot help us to make sense of how someone can learn to use those expressions.

This is why Fodor thinks that we need something more than truth conditions such as TC in order to make sense of linguistic comprehension and language learning. We need rules that will tell us which individual the name “Felicia” refers to, and which property is named by the predicate “  is tall.” If these rules are to be learnable then they must be stated in terms of expressions that the language user is already familiar with. In fact, we really need something like the following rule.

TC\* “Felicia is tall” is true just if X is G

Here “X” stands for another name for Felicia – one that the language user already understands (perhaps “X” might be “George’s sister”). Likewise “G” stands for another way of naming the property of being tall (perhaps “G” might be “greater than average in height”). This is what Fodor calls a truth rule.



### **Exercise 10.1** Explain in your own words the difference between the truth condition TC and the truth rule TC\*.

So, putting all this together, Fodor argues that learning a language has to involve learning truth rules. He thinks that this places some very fundamental constraints on any information-processing account of language learning. Learning a truth rule such as TC\* is, he thinks, a matter of forming hypotheses about what the expressions “Felicia” and “  is tall” mean. These hypotheses are then tested against further linguistic data and revised if necessary. Learning that George has no sisters, for example, would force me to revise my first version of the Felicia truth rule.

This is where the language of thought is required, Fodor argues. Learning a public language such as English, even if it is your first language, requires you to formulate, test, and revise hypotheses about the truth rules governing individual words. These hypotheses have to be formulated in some language. A truth rule is, after all, just a sentence. But which language are truth rules formulated in?

Fodor thinks that it cannot be the language being learned. You cannot use the language that you are learning to learn that language. That would be pulling yourself up by your own bootstraps. And since Fodor takes his account to apply to children learning their first language no less than to people learning a second language, the language cannot be any sort of public language. It can only be the language of thought, as described in Chapter 4.

That opens up the question: How is the language of thought itself learned? The short answer is that it cannot be learned. Since it is a language, the process of learning it would have to involve formulating hypotheses about truth rules (among other things). But Fodor is adamant that those hypotheses can only be formulated in the language of



thought. So – you would need to possess a language of thought in order to learn it. He draws the inevitable conclusion, which is that language of thought must be innate. We are all born with a fully developed language of thought, and it is this that allows us to learn the natural language of the community that we are born into.

Fodor's argument leads, therefore, to a form of *innatism* or *nativism* about language, because he thinks that the very possibility of language learning requires an innate language of thought. There are many other versions of innatism popular among linguists and cognitive scientists. All are based, in one form or another, on arguments to the effect that young children are simply not exposed to enough information, or even the right kind of information, to allow them to learn a language.

These arguments are collectively known as *poverty of the stimulus arguments*. The most famous proponent is the linguist Noam Chomsky, who uses them to support his general theory that all humans share a single language faculty, incorporating specialized language acquisition tools. Chomsky's view has gone through many iterations, but the basic idea has remained constant.

Chomsky thinks that all human languages (including sign language) can be understood in terms of different parameter settings in a universal grammar (the parameters are, as it were, optional settings, while the grammar provides a fixed structure that holds across all languages). What we think of as language learning is really parameter setting. The universal grammar is innate, and so all that a child needs to learn are the specific settings for their linguistic community. Chomsky views this process of parameter setting as involving processes of hypothesis formation and testing, rather similar to those proposed by Fodor.

This general picture of how language works is supported both by specific analyses of different languages and by poverty of the stimulus arguments. The details of Chomsky's different models of universal grammar are too complicated to go into here. But we can summarize the basic elements of his poverty of the stimulus arguments. He and his followers typically emphasize the following features of the young child's learning environment.

- Children are not positively rewarded for language learning.
- Children are typically only exposed to positive information (i.e., they are not told what counts as ungrammatical, but only given examples of grammatical utterances).
- The date from which each child learns is highly idiosyncratic.
- Much of the speech that children hear is actually ungrammatical, but not flagged as such.
- No child encounters more than a tiny fraction of the linguistic information that fixes the grammatical structure of the language.

It would be hard to dispute any of these claims about the learning environment for young children. But do they provide support for some version of innatism?

This seems to be an area where the proof of the pudding is in the eating. In other words, the best way to respond to an impossibility argument is by trying to give examples of exactly the sort of things that are being claimed to be impossible. In this case, these would be examples of how significant segments of language can be learned with little or no innate knowledge. We will be looking at examples in the remainder of this chapter.

 10.3

## Language Learning in Neural Networks

As we saw earlier, models of language learning are inextricably linked to models of linguistic understanding. Both Fodor and Chomsky have a very rules-focused way of thinking about linguistic understanding. This section introduces the very different connectionist approach to language mastery and language acquisition, which is based on modeling linguistic abilities in artificial neural networks.

When we first met artificial neural networks in Chapter 5, we discovered some fundamental differences between artificial neural networks and the sort of computational systems to which the physical symbol system hypothesis applies. In particular, we highlighted the following three differences.

- Representation in neural networks is distributed across the units and weights, whereas representations in physical symbol systems are encoded in discrete symbol structures.
- There are no clear distinctions in neural networks either between information storage and information processing or between rules and representations.
- Neural networks are capable of sophisticated forms of learning. This makes them very suitable for modeling how cognitive abilities are acquired and how they evolve.

The second feature in particular suggests that neural network models of language are going to be very different from the rules-based approach just discussed. If neural networks do not admit a clear distinction between rules and representations, then they cannot incorporate truth rules (or any other type of rules).

Despite this, neural networks have been strikingly successful at modeling language mastery, particularly when it comes to modeling how languages are learned (which is not surprising, of course, since learning is what neural networks are best at). In this section we will look at some influential and important studies. The networks in these studies show that there is an alternative to the rule-based conception of language comprehension and learning discussed in the previous section.

In fact, neural network models are widely held to have shown a system can reveal complex linguistic skills without having any explicit linguistic rules encoded in it. So, for example, the simple recurrent networks developed by Jeff Elman have been successfully trained to predict the next letter in a sequence of letters, or the next word in a sequence of words. This in itself is very important in thinking about the information processing involved in language learning. At the very least it casts doubt on claims that we can only think about language in terms of rule-based processing.

But researchers in this area have also made a second, very important, contribution. This contribution speaks more directly to issues about the psychological plausibility of neural network models. Developmental psychologists and psycholinguists have carefully studied patterns in how children learn languages. They have discovered that, in many aspects of language acquisition, children display a very typical trajectory. So, for example, children make very similar types of error at similar stages in learning particular grammatical constructions. Neural network researchers have explored the extent to which their models can



reproduce these characteristic patterns. They have found some striking analogies between how children learn and how neural networks learn. Let's look at them.

## The Challenge of Tense Learning

One of the most formidable problems confronting children learning a language such as English is that it has both regular and irregular verbs. Some verbs behave in very predictable ways. So, for example, their past tenses are formed according to straightforward rules. Consider the verb “to bat,” for example. This is a regular verb. In the present tense we have “I bat.” In the past tense this becomes “I batted.” There is a very simple rule here. For regular verbs we form the past tense by adding the suffix “-ed” to the stem of the verb. The stem of “to bat” is “batt-.” For regular verbs, then, all that one needs to know in order to be able to put them in the past tense is their stem.

Contrast regular verbs with irregular verbs. We have “I give” in the present tense. This becomes “I gave” in the past tense – not “I gived,” as the simple rule might suggest. Likewise for “I take,” which becomes “I took.” Irregular verbs, by their very nature, are not easily summarized by simple rules. It is true that there are observable regularities in how the past tenses of irregular verbs are formed. So, for example, we see that both “I ring” and “I sing” have similar past tenses (“I rang” and “I sang”). It would be unwise, however, to take this as a general rule for verbs ending in “-ing.” The past tense of “I bring” is most certainly not “I brang.” Anyone who has ever learned English as a second language will know that the corpus of irregular verbs is full of “false friends” such as these.

Yet somehow, more or less all young children in the English-speaking world manage to find their way through this minefield. How do they do it?

Researchers such as the psychologist Stan Kuczaj have studied the grammaticality judgments that children made about sentences involving past tense verbs in order to examine how their understanding of the past tense develops. The test sentences included both correct past tense forms (such as “brought” and “gave”) and incorrect ones (such as “brang” and “gived”). The incorrect ones were typically constructed either by treating irregular verbs as if they were regular (as in “gived”), or by exploiting “false friends” (as in “brang”). Looking at patterns of grammaticality judgments across populations of children aged from 3 to 11 has led researchers to hypothesize that children go through three distinct stages in learning the past tense.

In the first stage, young language learners employ a small number of very common words in the past tense (such as “got,” “gave,” “went,” “was”). Most of these verbs are irregular and the standard assumption is that children learn these past tenses by rote. Children at this stage are not capable of generalizing from the words that they have learned. As a consequence, they tend not to make too many mistakes. They can't do much, but what they do they do well.

In the second stage children use a much greater number of verbs in the past tense, some of which are irregular but most of which employ the regular past tense ending of “-ed” added to the root of the verb. During this stage they can generate a past tense for an invented word (such as “rick”) by adding “-ed” to its root. Surprisingly, children at this

**TABLE 10.1** The stages of past tense learning according to verb type

	<b>STAGE 1</b>	<b>STAGE 2</b>	<b>STAGE 3</b>
<b>Early verbs</b>	Correct	Overregularization errors	Correct
<b>Regular verbs</b>		Correct	Correct
<b>Irregular verbs</b>		Overregularization errors	Improvement with time

stage take a step backward. They make mistakes on the past tense of the irregular verbs that they had previously given correctly (saying, for example, “gived” where they had previously said “gave”). These errors are known as *overregularization* errors.

In the third stage, children cease to make these overregularization errors and regain their earlier performance on the common irregular verbs while at the same time improving their command of regular verbs. Table 10.1 shows the basic trajectory.

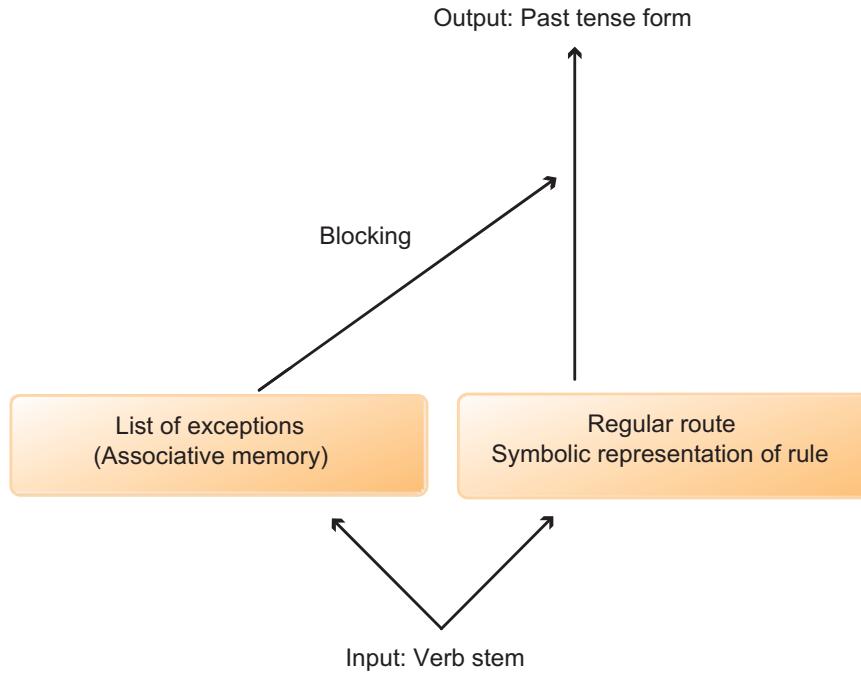
At first sight, this pattern of performance seems to support something like Fodor’s rule-governed conception of language learning. One might think, for example, that what happens in the second stage is that children make a general hypothesis to the effect that all verbs can be put in the past tense by adding the suffix “-ed” to the root. This hypothesis overrides the irregular past tense forms learned earlier by rote and produces the documented regularization errors. In the transition to the third stage, the general hypothesis is refined as children learn that there are verbs to which it does not apply and, correspondingly, begin to learn the specific rules associated with each of these irregular verbs.

The cognitive scientists Steven Pinker and Alan Prince have in fact proposed a model of understanding of the English past tense that fits very well with this analysis. Their model has two components and, correspondingly, two information-processing routes. These are illustrated in Figure 10.1.

One route goes via a symbolic representation of the rule that the past tense is formed by adding “-ed” to the stem of the verb. The symbolic component is not sensitive to the particular phonological form of the verb. It does not recruit information that, for example, the present tense of the verb ends in “-ing.” It simply applies the rule to whatever input it gets.

The second route, in contrast, goes via an associative memory system that is sensitive to the phonological form of the verb stem. It is responsible for storing exceptions to the general rule. It classifies and generalizes these exceptions in terms of their phonological similarity. One would expect this mechanism to pick up very quickly on the similarity, for example, between “sing” and “ring.”

The two routes are in competition with each other. The default setting, as it were, is the symbolic route. That is, the system’s “working hypothesis” is that it is dealing with a verb where the past tense is formed by adding “-ed” to the stem. But this default setting can be overridden by a strong enough signal coming from the associative memory system that keeps track of exceptions. What makes signals from the override system strong is that they have been suitably reinforced through experience. If I have had plenty of exposure to the



**Figure 10.1** The dual-route model of past tense learning in English proposed by Steven Pinker and Alan Prince.

“sing–sang” and “ring–rang” pairs, then this will strengthen the signal for “bring–brang.” But the more exposure I have to the “bring–brought” pair, the weaker the signal for “bring–brang.” Gradually, as I become increasingly exposed to different irregular forms, the signals that are reinforced end up being generally correct.

The model proposed by Pinker and Prince is certainly compatible with the general trajectory of how children learn the English past tense. It is also supported by the general considerations we looked at earlier. But should we accept it (or some other rule-based model like it)?

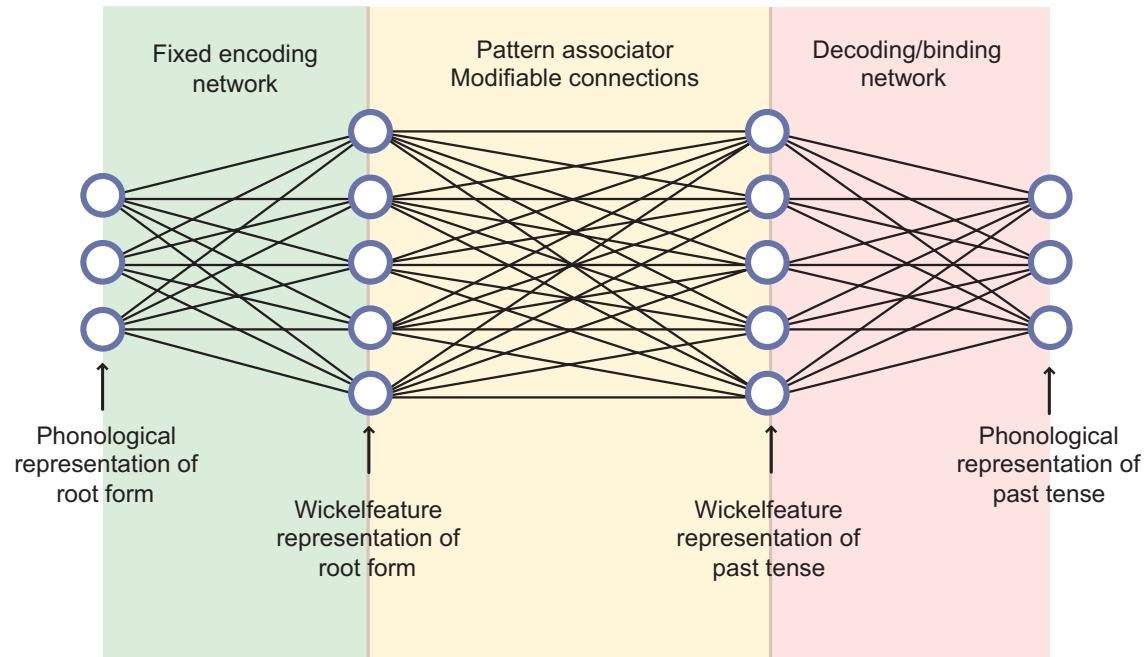


**Exercise 10.2** Explain how this two-component model of past tense understanding is compatible with the stages identified earlier in young children’s learning of the past tense in English.

This is where artificial neural networks come back into the picture, because researchers in neural network design have devoted considerable attention to designing networks that reproduce the characteristic pattern of errors in past tense acquisition without having programmed into them any explicit rules about how to form the past tense of verbs, whether regular or irregular.

## Connectionist Models of Tense Learning

The pioneering network in this area was designed by David Rumelhart and Jay McClelland and appeared in their 1986 collection of papers on parallel distributed processing. It was a relatively simple network, without any hidden units (and hence not requiring a



**Figure 10.2** Rumelhart and McClelland's model of past tense acquisition. (Adapted from Rumelhart, McClelland, and PDP Research Group 1986)

backpropagation learning algorithm), but nonetheless succeeded in reproducing significant aspects of the learning profile of young children. The network is illustrated in Figure 10.2.

There are really three different networks here. The first network takes as input a phonological representation of the root form of a verb. That is, it takes as input a sequence of phonemes. Phonemes are what linguists take to be the most basic meaningful constituents of words. An example is the phoneme /n/, which is the final sound in the words “tin” and “sin.” The first network translates this sequence of phonemes into a representational format that will allow the network to detect relevant similarities between it and other verb roots – as well as between the root forms and the correct past tense forms.

This representational format exploits an ingenious device that Rumelhart and McClelland call *Wickelfeatures* (after the cognitive psychologist Wayne Wickelgren, whose ideas they adapted). The details are very complex, but the basic idea is that a Wickelfeature representation codes phonetic information about individual phonemes within a word and their context. The aim is to represent verb stems in a way that can capture similarities in how they sound (and hence better represent the sort of stimuli to which young children are exposed).

The first network (the network converting phonological representations into Wickelfeature representations) is fixed. It does not change or learn in any way. The learning proper takes place in the second network. As the diagram shows, this network has no hidden units. It is a simple *pattern associator* mechanism. It associates input patterns with output patterns. The output patterns are also Wickelfeature representations of words, which are



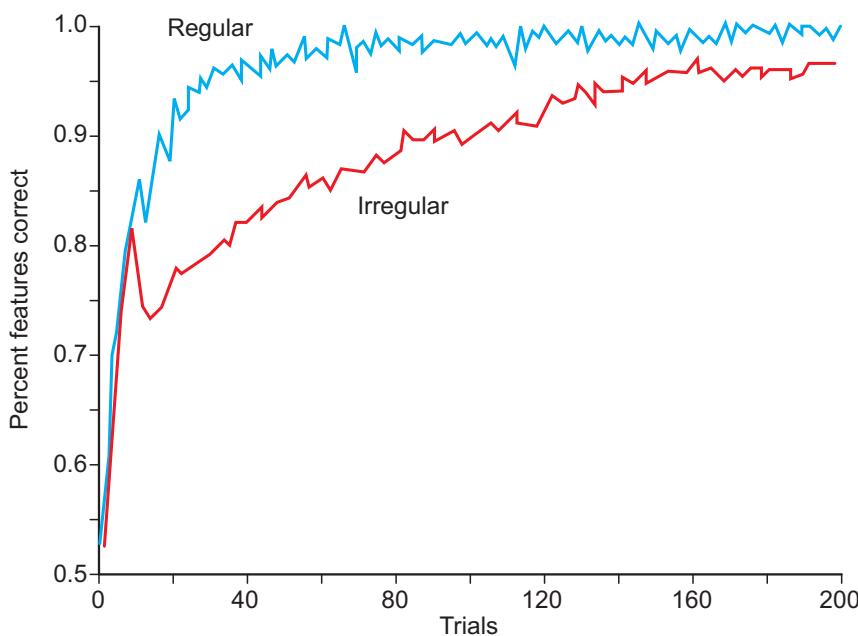
then decoded by the third network. This third network essentially reverses the work done by the first network. It translates the Wickelfeature representations back into sequences of phonemes.

The network was initially trained on ten high-frequency verbs. The aim here was to simulate the first stage in past tense acquisition. And then it was trained on a much larger training set of 410 medium-frequency verbs (of which 80 percent were regular).

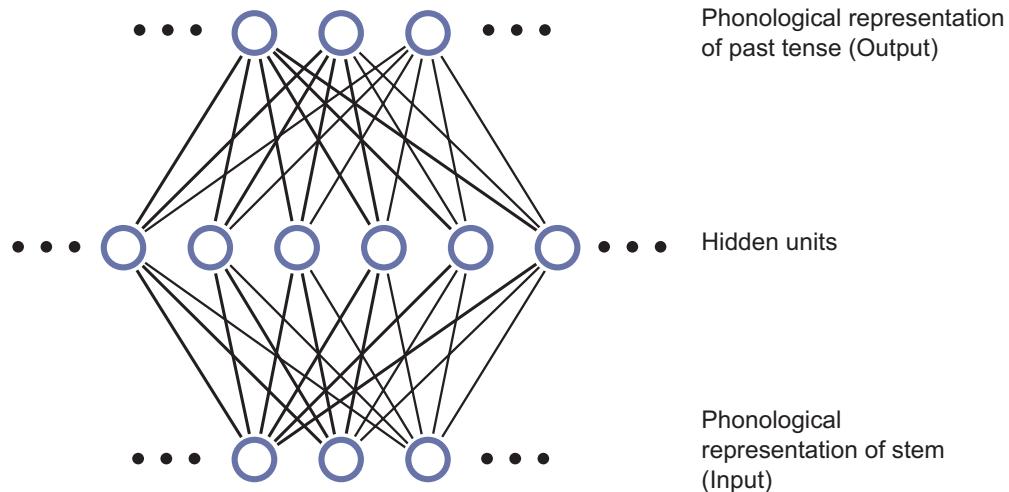
The learning algorithm used by the network is the perceptron convergence rule that we studied back in Section 5.2. At the end of the training the network was almost errorless on the 420 training verbs and generalized quite successfully to a further set of 86 low-frequency verbs that it had not previously encountered (although, as one might expect, the network performed better on novel regular verbs than on novel irregular verbs).

One interesting and important feature of the Rumelhart and McClelland network is that it reproduced the overregularization phenomenon. This is shown in Figure 10.3, which maps the network's relative success on regular and irregular verbs. As the graph shows, the network starts out rapidly learning both the regular and the irregular past tense forms. And then there is a sharp fall in performance on irregular verbs after the eleventh training cycle, while at the same time the degree of success on regular verbs continues to increase.

In other words, while the network's "understanding" of irregular verbs is "catching up" with its performance on regular verbs, it is still making characteristic errors. These errors involve treating irregular verbs as if they were regular. So, the network seems to be doing



**Figure 10.3** Performance data for Rumelhart and McClelland's model of past tense learning. The graph shows the success rates for both regular and irregular verbs. The line for irregular verbs clearly indicates the overregularization phenomenon. (Adapted from Rumelhart, McClelland, and PDP Research Group 1986)



**Figure 10.4** The network developed by Plunkett and Marchman to model children’s learning of the past tense. The network has a layer of thirty hidden units and is trained using the backpropagation learning algorithm. (Adapted from Plunkett and Marchman 1993)

exactly what young children do when they shift from the correct “gave” to the incorrect “gived” as the past tense of “give.”



**Exercise 10.3** Explain in your own words why it is significant that the Rumelhart and McClelland network produces the overregularization phenomenon.

Although the results produced by the Rumelhart and McClelland network are very striking, there are some methodological problems with the design of their study. In particular, as was pointed out in an early critique by Pinker and Prince, the overregularization effect seems to be built into the network. This is because the training set is so dramatically expanded after the tenth cycle. And since the expanded training set is predominantly made up of regular verbs, it has seemed to many that something like the overregularization phenomenon is inevitable.

Nonetheless, it is significant that a series of further studies have achieved similar results to Rumelhart and McClelland with less question-begging assumptions. Kim Plunkett and Virginia Marchman, for example, have produced a network with one layer of hidden units that generates a close match with the learning patterns of young children. The network is illustrated in Figure 10.4.

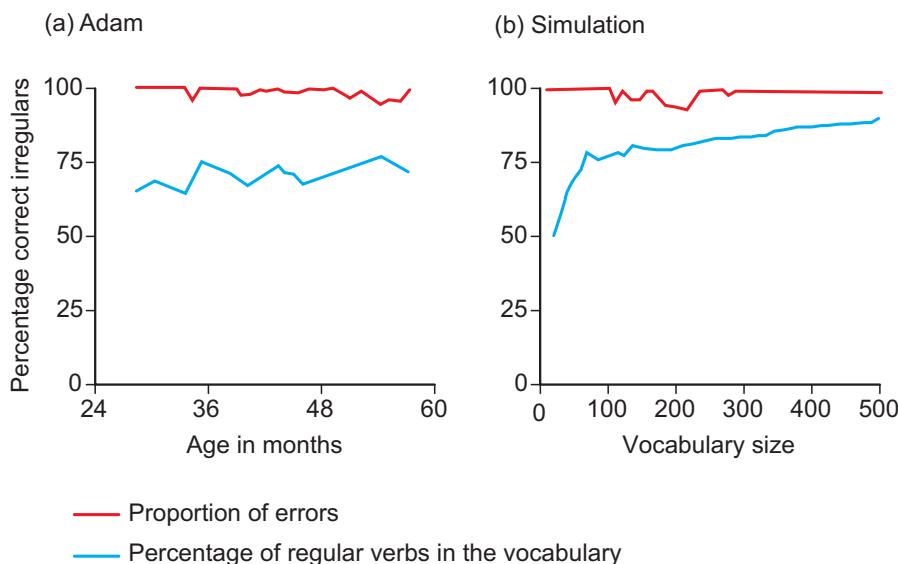
The Plunkett and Marchman network is in many ways much more typical of the type of neural network that are generally studied. Whereas the Rumelhart–McClelland network is a simple pattern associator using the perceptron convergence learning rule, the Plunkett–Marchman model has hidden units. Their model has twenty input and twenty output units. Between them is a single hidden unit layer with thirty units. The network uses the backpropagation learning algorithm. One advantage of this is that it removes the need to translate the initial phonological representation into Wickelfeatures.



Unlike the McClelland and Rumelhart model, the first stage of the training schedule was on twenty verbs, half regular and half irregular. After that the vocabulary size was gradually increased. There was no sudden increase – and hence no “predisposition” toward regularization errors. The percentage of regular verbs in the total vocabulary was 90 percent, which matches more or less the relative frequency of regular verbs in English. And yet the network did indeed display the characteristic trajectory, including the regularization errors characteristic of stage 2 learning in children. Plunkett and Marchman correctly guessed that the simple presence in the training set of both regular and irregular verbs would be enough to generate regularization errors during the second stage of training.

It is interesting to compare the learning profile of the Plunkett and Marchman network with the detailed profile of the learning pattern of a child studied by the psychologist Gary Marcus. The graph in Figure 10.5 compares the percentage of correctly produced irregular past tenses in the Plunkett and Marchman simulation and in a child whose past tense acquisition was studied by Marcus and colleagues. As the graph shows, the percentage of correctly produced irregular past tenses drops in both the network and the child as the vocabulary size increases. This seems to correspond to the second of the three stages identified earlier and to be correlated with the predominance of overregularization errors.

Certainly, there are huge differences between children learning languages and artificial neural networks learning to correlate verb stems with the correct versions of the past tense. And even when taken on their own terms, neural network models of language acquisition are deeply controversial, not least because of concerns about the biological plausibility of neural networks. But even with these caveats, using artificial neural networks to model



**Figure 10.5** A comparison of the errors made by Adam, a child studied by the psychologist Gary Marcus, and the Plunkett–Marchman neural network model of tense learning. Unlike the Rumelhart–McClelland model, this model uses hidden units and learns by backpropagation. (Adapted from McLeod, Plunkett, and Rolls 1998)

cognitive tasks offers a way of putting assumptions about how the mind works to the test – the assumption, for example, that the process of learning a language is a process of forming and evaluating hypotheses about linguistic rules.

The aim of neural network modeling is not to provide a model that faithfully reflects every aspect of neural functioning, but rather to explore alternatives to dominant conceptions of how the mind works. If, for example, we can devise artificial neural networks that reproduce certain aspects of the typical trajectory of language learning without having encoded into them explicit representations of linguistic rules, then that at the very least suggests that we cannot automatically assume that language learning is a matter of explicitly forming and testing hypotheses about linguistic rules. We should look at artificial neural networks not as attempts faithfully to reproduce the mechanics of cognition, but rather as tools for opening up new ways of thinking about how information processing might work.

## 10.4 Bayesian Language Learning

This final section turns to another alternative to explicit, rule-based models of language learning. This approach applies the Bayesian framework outlined in Chapter 7. As with the connectionist approach just considered, Bayesian models of language acquisition set out to show how the complexities of language can actually be learned without needing to postulate an innate language of thought or a universal language faculty.

According to Bayesians, language learning takes place through sensitivity to statistical regularities in heard speech, interpreted in terms of the different elements of Bayes's Rule – prior probabilities and likelihoods. The basic idea is that young children learning their first language (and adults learning a subsequent language) proceed by updating their probabilities according to Bayes's Rule – or, more accurately, that language learning can be modeled as a process of updating probabilities according to Bayes's Rule.

A quick refresher on Bayes's Rule. Bayes's Rule assumes that we have some evidence. In this case, the evidence is heard speech. It also assumes a hypothesis about how to interpret the evidence. The end result of applying Bayes's Rule is to derive what is called a *posterior probability*. The posterior probability is the probability of the hypothesis in the light of the evidence (or, more technically, the probability of the hypothesis conditional upon the evidence). It is called *posterior* because it comes *after* considering the evidence.

To apply Bayes's Rule, we need to have assigned a probability to the *likelihood* of the evidence conditional upon the hypothesis. That is, we need to have a view about how likely it is that we would see the evidence that we see if the hypothesis were true. We also need to have assigned a *prior* probability to the hypothesis. Our *prior* is the probability that we assign to the hypothesis *before* considering the evidence.

With all that in mind, Bayes's Rule says that the posterior probability of the hypothesis =

$$\frac{\text{Likelihood of the evidence} \times \text{Prior probability of the hypothesis}}{\text{Probability of the evidence}}$$



Bayes's Rule is a straightforward consequence of the definition of conditional probability.



### Exercise 10.4 Go back to Chapter 7 and review the discussion of Bayes's Rule.

Why might one think that Bayes's Rule would be a useful tool for thinking about language learning? A good starting point is studies showing that young children seem to be very sensitive to statistical regularities in heard speech, and indeed that adults can use statistical regularities to detect phrase structure.

## Probabilities in Word and Phrase Segmentation

One of the most basic challenges in making sense of speech is what is called *word segmentation*. Speech is a continuous stream of sound. In order to make sense of it, the continuous stream of sound needs to be segmented into individual words. This is obviously the first step in understanding language, an essential preliminary to decoding syntactic/grammatical structure and assigning meanings to words and sentences. And in fact, it emerges relatively early in infant development. Developmental linguists hold that word segmentation starts to emerge when infants are around 8 months old.

From a phonetic point of view, individual syllables are the natural breaks in the stream of speech. And they are also the building blocks of words. But how does an 8-month-old infant figure out which combinations of syllables make words, and which ones do not?

An influential model in developmental linguistics appeals to what are called *transitional probabilities*. Let's take three syllables /mo/, /ma/, and /pa/. The transitional probability between any two of them is the probability that the second will follow the first. In the language of probability, it is the probability of the second, conditional upon the first. So, for example, the transitional probability of /mo/ and /ma/ is the probability that /ma/ follows /mo/, while the transitional probability of /mo/ and /pa/ is the probability that /pa/ follows /mo/.

The basic idea behind applying transitional probabilities to word segmentation is that high transitional probabilities will tend to indicate syllables occurring within a word, while low transitional probabilities will tend to occur across the boundaries of words. This makes good sense, because there is much greater scope for variation between a syllable at the end of one word and a syllable at the beginning of the next than there is between two syllables occurring within a word.

But where do the probabilities come from, you may wonder? The answer is that they come from frequencies. Young infants have a lot of time to listen to the speech of those around them, sometimes directed at them but most often not. They also imitate and reproduce the sounds that they hear when they babble. As shown in the 1990s by Jenny Saffran, Richard Aslin, and Elissa Newport, infants, like all higher animals, are exquisitely sensitive to the frequency of correlations, and they can exploit this sensitivity to parse streams of sound into words.

In subsequent work, Elissa Newport and Susan Thompson looked at how transitional probabilities could be a factor in the much more complicated language learning problem of figuring out how to identify phrases. The basic idea is the same, except that now it is being applied to whole words, instead of to individual syllables. The transitional probabilities between two words will typically be higher when those words occur inside a phrase, than when they occur across the boundaries of a phrase. So, for example, take any noun phrase – “the camera” for example. The probability of the word “camera” coming after the word “the” will typically be much higher than the transitional probability of “the” coming after “camera.” This reflects the fact that nouns are typically used with articles (either definite, such as “the,” or indefinite, such as “a”) within a phrase, whereas the occurrence of “the” after “camera” would typically signal the start of a new phrase, or even a new sentence.

This work shows the importance of sensitivity to probability and frequencies in language learning. But it is not yet Bayesian reasoning, properly speaking. For that we need to turn to different examples.

## Understanding Pronouns

Linguists have devoted much time and energy to studying pronouns, such as “he,” “she,” “it,” and “one.” They have been particularly intrigued by the phenomenon typically called *pronominal anaphora* (and in fact, this is a case study often used in poverty of the stimulus arguments). A pronoun is used anaphorically when it picks up on the reference of a term earlier in the sentence. Here is a simple example (or rather, two examples): “Susan ate an ice cream and she enjoyed it.” In this sentence, the pronoun “she” refers to Susan and “it” refers to the ice cream that Susan ate.

It is not difficult to figure out what the two pronouns refer to in “Susan ate an ice cream and she enjoyed it.” Once you’ve noticed that “she” is a feminine pronoun, everything falls into place very quickly. But many examples of pronominal anaphora are not so straightforward. Consider these two, for example.

- (1) Tommy beat Mike, and he didn’t like it.
- (2) Susan has a black cat, and I want one.

There is no problem interpreting the pronoun “it” in (1). It refers to the fact that Tommy beat Mike. But the second pronoun “he” can be interpreted in two different ways. It could refer either to “Tommy” or to “Mike.” So, as written, (1) is ambiguous between two possible references for “it.”

Sentence (2) is ambiguous in a different way. What is it that I want? I don’t want Susan’s exact cat. But do I specifically want a black cat, like Susan’s, or am I just generally cat-deprived, so that any old cat will do? You can think about it like this. (2) is ambiguous between two possible membership classes for the pronoun “one.” The pronoun “one” could refer to a member of the class of cats, or to a member of the more specific class of black cats.



Typically, though, hearers manage to resolve both ambiguities. How do they do it? This problem has been tackled from a Bayesian perspective. Here is another sentence with anaphoric ambiguity.

- (3) I'll play with this red ball and you can play with that one.

The pronoun “one” is ambiguous in (3) in exactly the same way as it was in (2). Are you being instructed to play with a red ball specifically, or with any old ball, which might or might not be red? The sentence still seems ambiguous, even if there is only one other ball in the vicinity, and that ball happens to be red. Even if you start playing with the red ball, that in itself doesn't settle whether what I meant was “You can play with that ball” or “You can play with that red ball.”

It turns out that this exact example was proposed by Jeffrey Lidz and collaborators as evidence that infants are born with innate knowledge of syntactic structure (in an article provocatively titled “What children know about syntax but could not have learnt”). According to Lidz, it is obvious that the second reading is correct (with the relevant membership class being red balls, rather than balls in general). But, he thinks, no evidence available to language users could possibly show that to be the intended reference, for the simple reason that all red balls are balls. This means that the evidence will always confirm both hypotheses. Yet, young children are able to disambiguate this type of sentence from a very early age, and so, he concludes, they must have used innate knowledge to do so.

Responding to Lidz et al., Terry Regier and Susan Gahl showed how Bayesian reasoning could solve the problem. They accept the basic argument that there can be no *direct* evidence to support the red ball hypothesis, but they think that Bayesian children might have other, *indirect* tools at their disposal. In particular, they point out that, from a Bayesian perspective, we need to look at the likelihoods. In other words, we need to consider how probable the evidence is, given the hypothesis.

Imagine that you hear someone say to you “I'll play with this red ball and you can play with that one.” You've heard that same sentence a few times before, and every time there has been a red ball there for you to play with. The only evidence you have, therefore, is that the ball you have to play with is red, and as we've said, that is perfectly compatible with the phrase “that one” referring back to the general class of balls, rather than to the specific class of red balls.

But you need to think about how likely the evidence is, given the hypothesis. In other words, you need to ask yourself which of these two scenarios is more probable:

- (a) The balls in the room are always red (*evidence*), if what he means is that I should play with that red ball (*hypothesis*)?
- (b) The balls in the room are always red (*evidence*), if what he means is that I should play with any old ball, red or not (*hypothesis*)?

Regier and Gahl make the plausible claim that, from a Bayesian perspective, (a) is more probable than (b). This is because it fits the evidence better. If the instruction is to play with

any old ball, irrespective of color, then you would expect to see some nonred balls in the room at least some of the time.

Here's the analogy they give. Suppose you are trying to decide whether all animals bark (*hypothesis 1*), or whether only dogs bark (*hypothesis 2*). Your evidence is that you have heard lots of dogs barking, but you have never heard anything that is not a dog bark. All dogs are animals, of course, and so your evidence is consistent with both hypotheses. But still, you would be wise to opt for hypothesis 2, because if all animals barked then you would expect to have run into some barking animals that are not dogs.

One of the key claims of Bayesian approaches to language learning, therefore, is that cognitive scientists have been too quick to make claims about aspects of language being impossible to learn. Children and adults can learn a lot from thinking about Bayesian likelihoods, as the example of pronominal anaphora shows.

## Learning Linguistic Categories

A related problem much studied by linguists and cognitive scientists is how children come to learn linguistic categories, and so distinguish between different types of noun. This can be a particularly challenging problem when the categories have overlapping membership. To continue with our canine example, all Dalmatians are dogs, and all dogs are animals. So, we have three linguistic categories:

Dalmation (subordinate category)

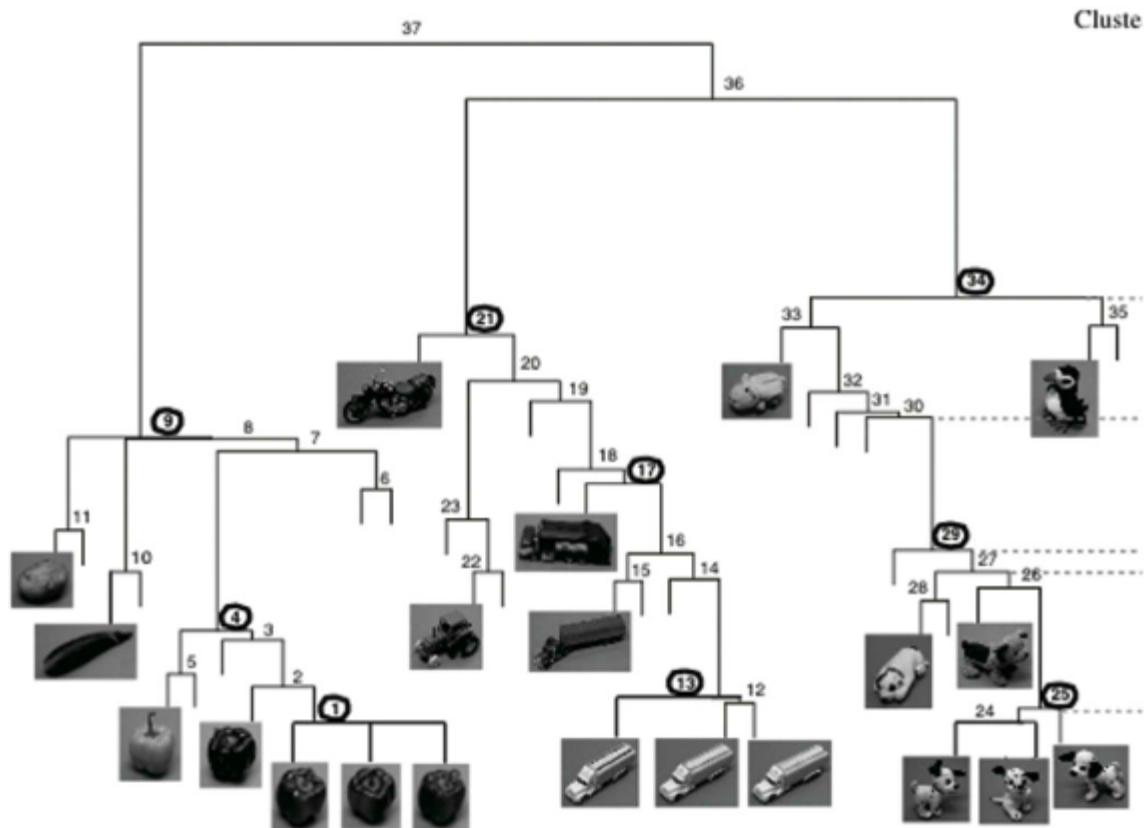
Dog (basic category)

Animal (superordinate category)

Imagine a child trying to learn the word "dog." The family dog happens to be Spotty the Dalmatian, and she learns to associate the word "dog" with Spotty. Since she knows that Spotty's name is "Spotty," she can quickly figure out that "dog" is not another name for Spotty. But how does she figure out that "dog" refers to all dogs, not just to Dalmatians – and not to all animals either?

This is another area where cognitive scientists have appealed to innate knowledge, suggesting that children are born with a bias in favor of interpreting category words as basic-level categories, corresponding for example to biological species. As before, this has been supported by a general argument that children could not possibly learn to discriminate between subordinate, basic, and superordinate categories simply on the basis of experience. This argument has been challenged by Fei Xu and Joshua Tenenbaum, among others.

Xu and Tenenbaum ran a set of experiments on category learning in adults to try to uncover tools that young children might be using to learn category words. The experimenters used photographs of objects that fell into subordinate-basic-superordinate



**Figure 10.6** A hierarchical cluster of similarity judgments, with nodes corresponding to clusters of stimuli more similar on average to each other than to objects in the nearest cluster.

taxonomies. One of these taxonomies was in fact the Dalmatian-dog-animal taxonomy (although the name of the Dalmatian was not recorded). They were asked to figure out the meaning of a new, non-English word (e.g., “blick”), based on a small set of labeled samples. This part of the experiment was intended to map the experience of a typical language-learning child.

At the same time, in an effort to understand the tools language users might be using, Xu and Tenenbaum asked their subjects to make pairwise similarity judgments (on a scale of 1 to 9) between the forty-five different objects in the photographs, instructing them to focus on the aspects of the objects that they had used in interpreting the new words. They then used these similarity judgments to draw a hierarchical cluster map of perceived similarities. An example is illustrated in Figure 10.6. Distance in the map corresponds to perceived similarity. Each node in the tree represents a cluster of stimuli more similar to each other on average than to objects in the nearest cluster.

Xu and Tenenbaum then developed a Bayesian model, using this cluster analysis as the hypothesis space. Each similarity cluster corresponded to a single hypothesis about the possible extension of a word such as “blick,” on the principle that one can’t view objects as

falling under a single category word without seeing them as similar in certain respects. This means that any set of objects viewed as similar is a candidate for the extension of a given word.

But how, given a range of hypotheses (corresponding to clusters of objects perceived as similar), is one going to be decided upon? Xu and Tenenbaum propose a Bayesian decision process, exploiting the key Bayesian concepts of *prior probability* and *likelihood*. They used the cluster space to derive priors and likelihoods.

The basic idea behind how they derived the priors is that the more distinctive a similarity cluster is, the more likely it is to have its own name. This makes good sense. Languages contain names to pick out salient groups of things, after all, and the more distinctive a similarity cluster is, the more it will stand out (be salient) – and the easier it will be for language users to pick out members of the group.

The likelihoods for the model were arrived at by applying what they called the *size principle*. This is in effect the same principle that we looked at in the context of red balls and barking dogs. It basically says that you should interpret the hypothesis being confirmed as narrowly as possible.

So, the hypothesis that “blick” refers to Dalmatians is far more likely than the hypothesis that it refers to dogs, if you have heard it applied to three Dalmatians in a row – and much more likely still than the hypothesis that “blick” refers to animals in general. Remember that the likelihood measures the probability that a hypothesis is true if the evidence is relevant. And what the size principle says is that a narrow hypothesis (relative to the evidence) is more likely to generate that evidence than a broader one.

The key finding from Xu and Tenenbaum’s experiments on adults was that their category membership judgments were accurately predicted by a Bayesian model using the priors and likelihoods just described – i.e., with priors fixed by the similarity space and likelihoods by the size principle. Moreover, when they went on to test 3- to 4-year-old children on a similar task, they found that the children’s categorization judgments closely mapped those of the adults, which is certainly consistent with the hypothesis that their linguistic behavior could be captured by a Bayesian model (assuming that their similarity judgments mapped those of adults).

The Xu and Tenenbaum experiments provide another illustration of how Bayesian statistical inference might be deployed by language learners to solve problems that many cognitive scientists have taken to be unsolvable by any type of learning.

This brings us back to the cluster of issues with which we began. How one thinks about language learning is very closely tied to how one thinks about linguistic understanding. At one end of the spectrum, represented by Jerry Fodor and Noam Chomsky, linguistic understanding is fundamentally a matter of mastering and applying linguistic rules. This position is often accompanied by poverty of the stimulus arguments to the effect that substantial parts of language learning must be innate, as young children do not encounter the sort of evidence that would be required to learn them. At the other end of the spectrum, the emphasis is less on explicit linguistic rules and more on learning mechanisms, such as neural networks or Bayesian principles.



## Summary

This chapter has explored different ways of modeling language mastery and language learning. We started out with a very general argument, due to Jerry Fodor, that language learning must be rule-based, and saw how that line of thinking led him to a form of innatism about language. Chomsky reaches a similar view about language learning, from a rather different starting point. We then turned to two very different approaches. With respect to the very specific problem of how children learn the past tense of English verbs, we saw how connectionist models of tense learning offer an alternative to the idea that grammar is learned by internalizing explicitly represented grammatical rules. The final section looked at Bayesian models of language learning, and showed how a Bayesian approach can illuminate word segmentation, pronominal anaphora, and learning hierarchically organized category words.

## Checklist

### Language and Rules

- (1) Language is a paradigmatically rule-governed activity (not just grammatical rules, but also rules giving the meanings of individual words and governing the deep structure of sentences).
- (2) The default hypothesis in thinking about language learning is that it is a matter of learning the rules that govern the meanings of words and how they combine into meaningful units.
- (3) Fodor has built on the default hypothesis to argue that learning a language requires learning *truth rules*, which must be stated in the language of thought.
- (4) According to Fodor, the language of thought cannot itself be learned, and so must be innate.
- (5) Noam Chomsky has reached a similar nativist/innatist conclusion based on poverty of the stimulus arguments.
- (6) One way to challenge such arguments is to construct models that simulate the trajectory of human language learning without explicitly representing any rules.

### Modeling the Acquisition of the English Past Tense

- (1) Children learning the English past tense go through three easily identifiable stages:  
Stage 1 They employ a small number of verbs with (mainly irregular) past tenses.  
Stage 2 They employ many more verbs, tending to construct the past tense through the standard stem + -ed construction (including verbs they had formerly got right).  
Stage 3 They learn more verbs and correct their overregularization errors.
- (2) This pattern of past tense acquisition can be accommodated by a symbolic model.
- (3) But connectionist models of past tense acquisition have been developed that display a similar trajectory without having any rules explicitly coded in them.

### Bayesian Language Learning

- (1) The basic Bayesian idea is that language learning can be modeled as a processing of updating probabilities according to Bayes's Rule.

- (2) The Bayesian approach is supported by studies showing that children and adults are both highly sensitive to statistical regularities in heard speech.
- (3) Studies suggest that infants use transitional probabilities for word segmentation (parsing heard speech into words), and that transitional probabilities are also used by adults to map the boundaries of phrases.
- (4) Pronominal anaphora (identifying the antecedent of a pronoun) has been proposed as a particular challenge for noninnatist approaches to language learning, but the problem seems susceptible to a broadly Bayesian solution.
- (5) Another traditional problem for empiricist theories of language learning is how children learn hierarchical category words that pick out overlapping sets of objects. Xu and Tenenbaum developed a Bayesian model to solve this problem with appropriate priors and likelihoods, showing that it captured the linguistic behavior of adults and 3- to 4-year-old children.

## Further Reading

Fodor's discussion of truth rules is in his book *The Language of Thought* (Fodor 1975). Chomsky first raised poverty of the stimulus arguments in his influential review of B. F. Skinner's book *Verbal Behavior* in 1959. See also Chomsky 1968, 1980a (summarized with comments and critique in Chomsky 1980b) and, for more recent discussion, Berwick et al. 2011. For a critical discussion of these arguments, see Pullum and Scholz 2002, and for a general discussion of arguments for nativism, see Cowie 1999.

The second volume of *Parallel Distributed Processing* (McClelland, Rumelhart, and the PDP Research Group 1986) contains a number of papers applying the theoretical framework of connectionism to different cognitive abilities. Some of these applications are explored further in McLeod, Plunkett, and Rolls 1998 and Plunkett and Elman 1997. For more general discussion of modeling within a connectionist framework, see Dawson 2004. Paul Churchland has been a tireless proponent of the power of connectionist networks; see, for example, the papers in Churchland 2007 for a wide range of applications. See also McClelland et al. 2010.

Chapter 18 of the original PDP collection (Rumelhart and McClelland 1986) was the first salvo in what has become a lengthy debate about how to model past tense learning. Pinker and Prince 1988a made some telling criticisms of Rumelhart and McClelland's model (Pinker and Prince 1988b, reprinted in Cummins and Cummins 2000, is more condensed). A number of researchers took up Pinker and Prince's challenge – see, for example, Plunkett and Marchman 1993. The work by Marcus described in the text is presented in Marcus et al. 1992. For a more recent exchange, see Pinker and Ullman 2002 and the reply in McClelland and Patterson 2002. Connectionist models have been applied to many different aspects of language. Plaut, Banich, and Mack 2003 describes applications to phonology, morphology, and syntax. Christiansen and Chater 2001 is an interdisciplinary collection of papers in the emerging field of connectionist psycholinguistics. Westermann and Ruh 2012 provides a review of different approaches to past tense learning, including connectionist approaches. Perhaps the most famous formal result in the theory of language learning is Gold's theorem, which places constraints upon the class of languages that can be learned with purely positive feedback. Gold's theorem is clearly presented in Johnson 2004.



Doug Rohde and David Plaut have used neural network models to argue that Gold's theorem cannot straightforwardly be applied in cognitive science (Rohde and Plaut 1999).

For a helpful and up-to-date overview of Bayesian approaches to language acquisition, see Pearl and Goldwater 2016. See also Chater and Manning 2006. Saffran, Aslin, and Newport 1996 and Aslin, Saffran, and Newport 1998 were pioneering studies of word segmentation in infancy, mainly working with artificial languages. Pelucchi, Hay, and Saffran 2009 extended the approach to more realistic linguistic situations. Thomson and Newport 2007 apply transitional probabilities to phrase segmentation. The Bayesian analysis of pronominal anaphora described in the text originates in Regier and Gahl 2004, who were responding to Lidz, Waxman, and Freeman 2003. The Xu and Tannenbaum experiments on category learning are described in Xu and Tannenbaum 2007. For a recent extension of this approach to one-shot category learning, see Lake, Salakhutdinov, and Tannenbaum 2016.





## CHAPTER ELEVEN

# Object Perception and Folk Physics

### OVERVIEW 285

- 11.1 Object Permanence and Physical Reasoning in Infancy** 286  
Infant Cognition and the Dishabituation Paradigm 286  
How Should the Dishabituation Experiments Be Interpreted? 292

- 11.2 Neural Network Models of Children's Physical Reasoning** 293  
Modeling Object Permanence 295  
Modeling the Balance Beam Problem 297

- 11.3 Conclusion: The Question of Levels** 300



## Overview

Back in Chapter 5, we saw how information processing works in single-unit networks and then looked at how the power of neural networks increases when hidden units are added. In Chapter 10 we started exploring how neural networks can model cognition. We looked at neural network models of past tense learning and saw how their learning trajectory bears striking resemblances to the learning trajectory of human infants. This chapter turns to another application of neural networks. We will see how they can be used to model object perception (and, in particular, what developmental psychologists call object permanence).

Many studies have shown that the perceptual universe of human infants is far more complex and sophisticated than was traditionally thought. From a very early age human infants seem to be sensitive to certain basic properties of physical objects. They have definite (and often accurate) expectations about how objects behave and interact. Some of this research is presented in Section 11.1, where we see how it can very naturally be interpreted in computational terms, as involving an explicitly represented and quasi-theoretical body of rules and principles (*a folk physics*).

In Section 11.2, however, we show how some of the very same data can be accommodated without this type of explicit, symbolic representation. We look at some neural network models that share some of the basic behaviors of the infants in the experiments without having any rules or

principles explicitly coded into them. This opens the door to a different way of thinking about infants' knowledge of the physical world.

Finally, in Section 11.3 we turn to an issue that has been in the background throughout this chapter and Chapter 10. How exactly should we think about the relation between symbolic models (physical symbol systems) and neural network models? Are they in competition with each other. Or are they giving accounts on different levels? We explore these issues through a famous objection to neural network modeling due to Jerry Fodor and Zenon W. Pylyshyn.

## 11.1

### Object Permanence and Physical Reasoning in Infancy

What is it like to be a human infant? Until fairly recently, most developmental psychologists were convinced that the infant experience of the world is fundamentally different from our own. The famous psychologist and philosopher William James (brother of the novelist Henry James) coined the memorable phrase “a blooming, buzzing, confusion” to describe what it is like to be a newborn infant (a *neonate*, in the jargon of developmental psychologists). According to James, neonates inhabit a universe radically unlike our own, composed solely of sensations, with no sense of differentiation between self and objects or between self and other, and in which the infant is capable only of reflex actions. It takes a long time for this primitive form of existence to become the familiar world of people and objects and for reflexes to be replaced by proper motor behavior.

The most famous theory within the traditional view was developed by the Swiss psychologist Jean Piaget (1896–1980). According to Piaget, infants are born with certain innate, reflex-like sensorimotor schemas that allow them to perform very basic acts such as sucking a nipple. Infants gradually bootstrap these basic schemas into more complex behaviors (what Piaget called circular reactions) and gradually come to learn that they inhabit a world containing other objects and other individuals. According to Piaget, infants are born highly egocentric and it is not until the end of what he called the sensorimotor stage (at around 2 years of age) that they come fully to appreciate the distinctions between self and other and between the body and other physical objects.

In recent years, however, researchers have developed new techniques for studying the cognitive abilities of neonates and older infants. These techniques have led to a radical revision of the traditional view. As a consequence, many developmental psychologists now think that the world of the human infant is much less of a “blooming, buzzing, confusion” than James thought. It is now widely held that even very young infants inhabit a highly structured and orderly perceptual universe. The most famous technique in this area is called the *dishabituation paradigm*, which is a technique for exploring the expectations that infants have about how objects will behave.

### Infant Cognition and the Dishabituation Paradigm

The basic idea behind the dishabituation paradigm is that infants look longer at events that they find surprising. So, by measuring the amount of time that infants look at events of



different types, experimenters can work out which events the infants find surprising and then use this to work backward to the expectations that the infants had about how those events were going to turn out.

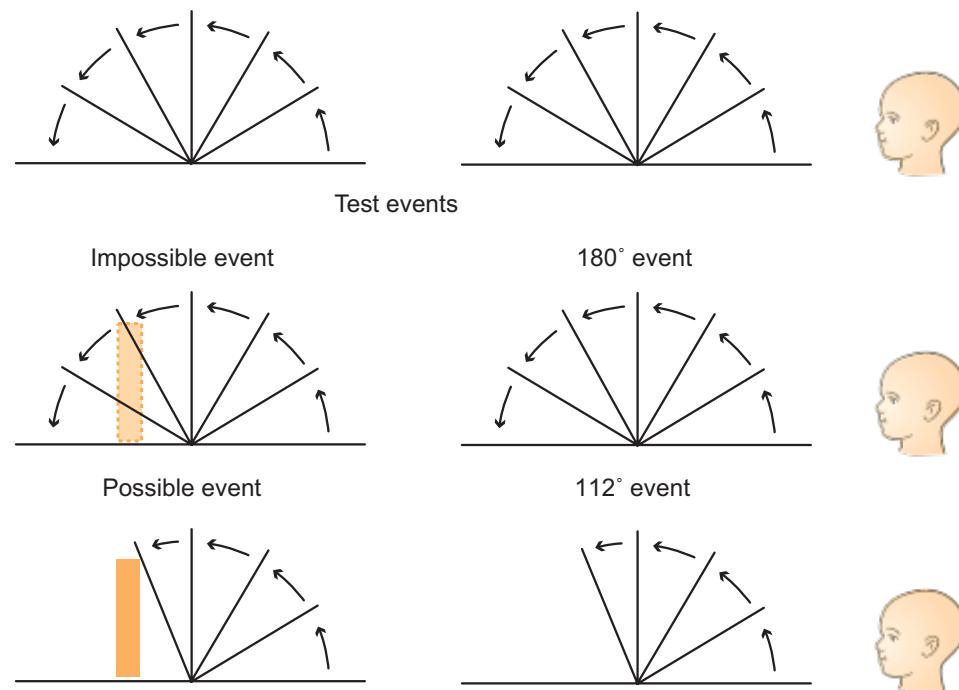
So, the basic idea is that if infants look longer at something, then that suggests that it did not turn out the way they expected. This basic idea is applied in practice in a number of ways. One technique is to habituate infants to a given type of event (i.e., presenting the infants with examples until they lose interest) and then to present them with events that differ from the original one in certain specified ways. Looking-time measures can then be used to identify which of the new events capture the infants' attention, as measured by the amount of time the infants spend looking at them. This allows experimenters to detect which features of the events the infants find surprising – and hence to work out how the infants expected the events to unfold. This way of identifying “violation of expectations” is called the dishabituation paradigm.

The developmental psychologist Renée Baillargeon devised a very influential set of experiments using the dishabituation paradigm. We can use her *drawbridge* experiments to illustrate how the paradigm works and what we can learn from it about the perceptual universe of the human infant. In one set of experiments, Baillargeon habituated her infants (who were all about 4.5 months old) to a screen (the drawbridge) rotating 180 degrees on a table. She was interested in how the infants would react when an object was hidden within the drawbridge's range of motion, since this would be a way of finding out whether the infant had any expectations about objects it could not directly perceive.

In order to investigate this, Baillargeon contrived a way of concealing the object so that, although it could not be seen by the infant, any adult or older child looking at the apparatus could easily work out that it would obstruct the movement of the screen. She then presented infants with two different scenarios. In the first scenario the screen rotated as it had done before until it got to the place where the obstructing box would be – and then it stopped, exactly as you or I would expect it to. In the second scenario, the screen kept on rotating for the full 180 degrees and hence apparently passed through the obstructing box. The experiments are illustrated in Figure 11.1.

Baillargeon found that the infants looked significantly longer in the second scenario. They were, it seemed, surprised that the screen looked as if it was passing straight through the obstructing box. So, if we assume that infants look longer when their expectations are violated, the experiments show that they do not expect the screen to keep on rotating through the place where the obstructing box would be. Baillargeon concluded that, although the infants could not see the obstructing box, in some sense they nonetheless “knew” that the box was there – and that the screen could not pass through it.

This result is very interesting because it has direct implications for a long-running debate in developmental psychology. Developmental psychologists have long been concerned with the question: At what stage, in early childhood or infancy, is it appropriate to ascribe a grasp that objects exist even when not being perceived? (Or, as developmental psychologists often put it, at what stage in development does *object permanence* emerge?) On the traditional view, derived ultimately from Piaget, object permanence does not appear until relatively late in development, at about 8 or 9 months. What Baillargeon's drawbridge



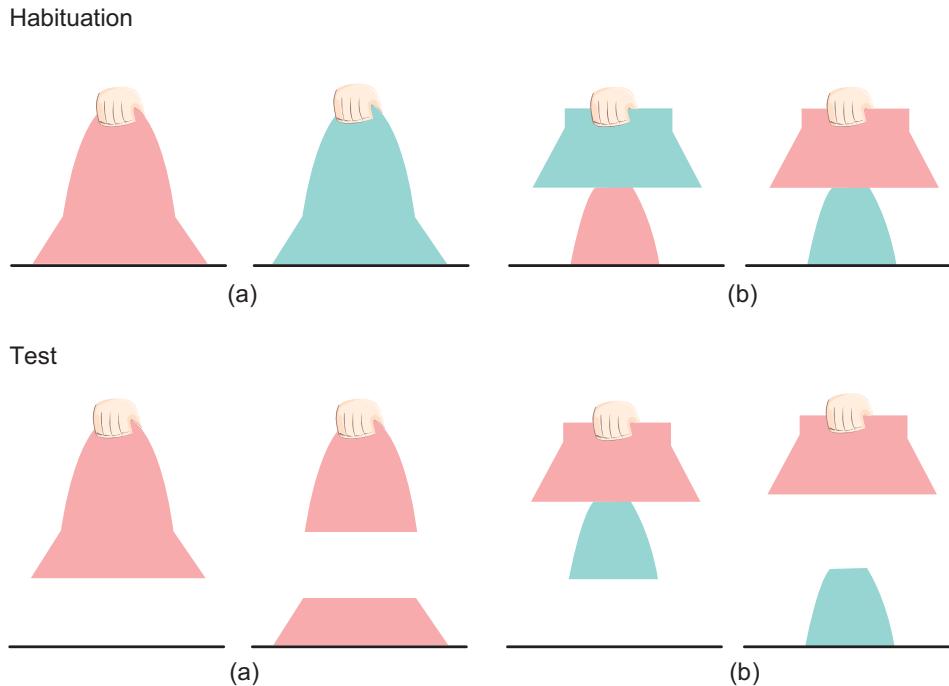
**Figure 11.1** Schematic representation of the habituation and test conditions in Baillargeon's drawbridge experiments. After habituation to a drawbridge moving normally through 180 degrees, infants were tested both on an impossible event (in which the drawbridge's movement would require it to pass through a hidden object) and a normal event (in which the drawbridge halts at the point where it would make contact with the hidden object). Baillargeon found that 4.5-month-old infants reliably looked longer in the impossible condition. (Adapted from Baillargeon 1987)

experiments seem to show, however, is that object permanence emerges much earlier than Piaget (and others) had thought.

But there is more going on here than simply object permanence. After all, it is not just that the infants are in some sense aware that the obstructing box is there even though they cannot see it. Their surprise at the second scenario shows that they have expectations about how objects behave. And, in particular, about how objects should interact. In fact, Baillargeon's drawbridge experiments, together with other experiments using the same paradigm, have been taken to show that even very young infants have the beginnings of what is sometimes called *folk physics* (or *naïve physics*) – that is to say, an understanding of some of the basic principles governing how physical objects behave and how they interact.

Elizabeth Spelke is another pioneer in using dishabituation experiments to study the perceptual universe of human infants. She has used a battery of experiments to argue that from a very young age infants are able to parse the visual array into spatially extended and bounded individuals. These individuals behave according to certain basic principles of physical reasoning. She thinks that four of these principles are particularly important for understanding the infant's folk physics.

The first principle is the *principle of cohesion*, according to which surfaces belong to a single individual if and only if they are in contact. It is evidence for the principle of



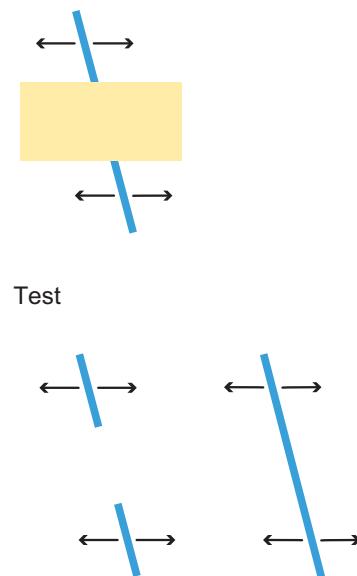
**Figure 11.2** Schematic representation of an experiment used to test infants' understanding of object boundaries and sensitivity to what Spelke calls the principle of cohesion (that surfaces lie on a single object if they are in contact). (Adapted from Spelke and Van de Walle 1993)

cohesion, for example, that infants do not appear to perceive the boundary between two objects that are stationary and adjacent, even when the objects differ in color, shape, and texture. Figure 11.2 illustrates how sensitivity to the principle of cohesion might be experimentally tested.

Three-month-old infants are habituated to two objects, one more or less naturally shaped and homogeneously colored, and the other a gerrymandered object that looks rather like a lampshade. When the experimenter picks up the objects, they either come apart or rise up cleanly. Infants show more surprise when the object comes apart, even if (as in the case of the lampshade) the object does not have the Gestalt properties of homogeneous color and figural simplicity. The conclusion drawn by Spelke and other researchers is that the infants perceive even the gerrymandered object as a single individual because its surfaces are in contact.

The principle of cohesion suggests that infants will perceive objects with an occluded center as two distinct individuals, since they cannot see any connection between the two parts. And this indeed is what they do – at least when dealing with objects that are stationary. Thus, it seems that infants do not perceive an occluded figure as a single individual, *if the display is static*. After habituation to the occluded figure they showed no preference for either of the test displays.

On the other hand, however, infants do seem to perceive a center-occluded object as a single individual if the object is in motion (irrespective, by the way, of whether the motion



**Figure 11.3** Schematic representation of an experiment testing infants' understanding of the principle of contact (that only surfaces in contact can move together). (Adapted from Spelke and Van de Walle 1993)

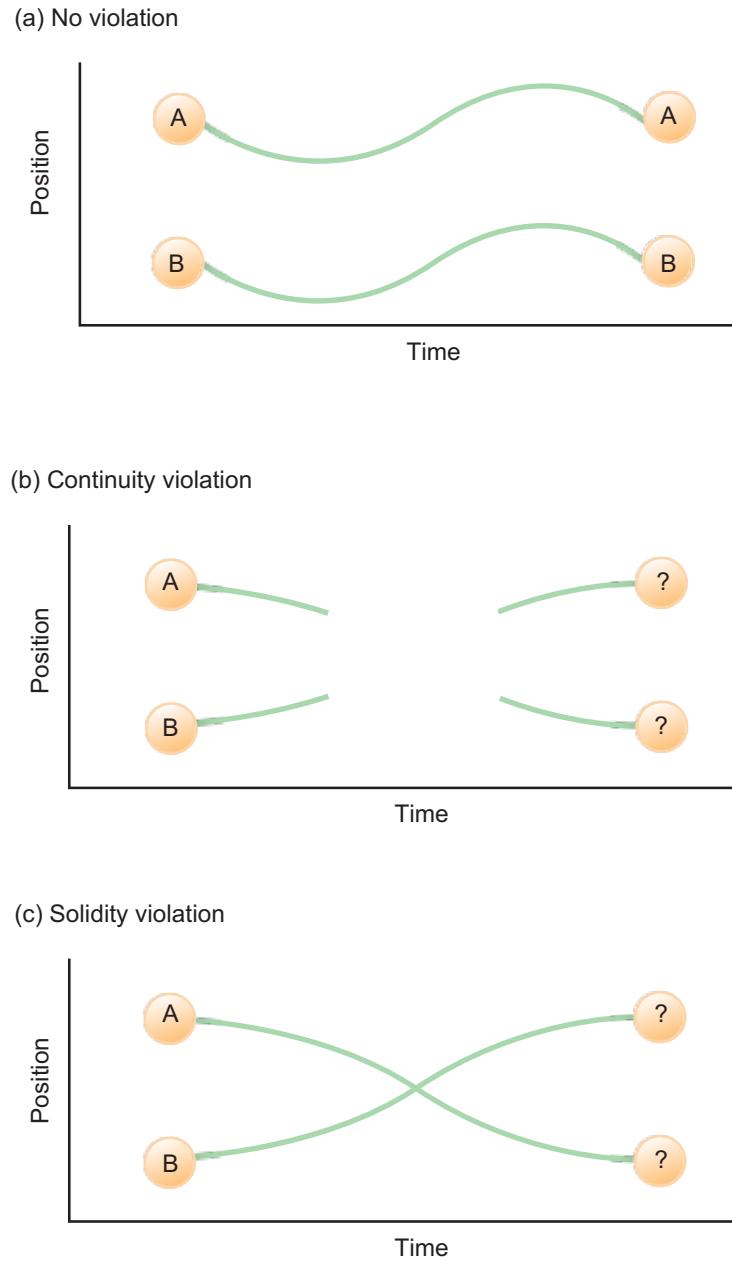
is lateral, vertical, or in depth). According to Spelke this is because there is another principle at work, which she terms the *principle of contact*. According to the principle of contact, only surfaces that are in contact can move together. When the principle of cohesion and the principle of contact are taken together they suggest that, since the two parts of the occluded object move together, they must be in contact and hence in fact be parts of one individual. This is illustrated in Figure 11.3.



**Exercise 11.1** Explain how an infant who understands the principles of cohesion and contact might respond to the two test situations depicted in Figure 11.3.

Spelke identifies two further constraints governing how infants parse the visual array. A distinctive and identifying feature of physical objects is that every object moves on a single trajectory through space and time, and it is impossible for these paths to intersect in a way that would allow more than one object to be in one place at a time. One might test whether infants are perceptually sensitive to these features by investigating whether they are surprised by breaches of what Spelke calls the *solidity* and *continuity* constraints. The drawbridge experiment that we have just discussed is a good example of reasoning according to the solidity constraint, since it shows that infants are sensitive to the impossibility of there being more than one object in a single place at one time. Figure 11.4 is a schematic representation of an experiment to test whether infants parse their visual array in accordance with the continuity and solidity constraints.

These are just some of the experiments that have been taken to show that even very young infants have a surprisingly sophisticated understanding of the physical world.



**Figure 11.4** Schematic depiction of events that accord with, or violate, the continuity or solidity constraints. Solid lines indicate each object's path of motion, expressed as changes in its position over time. Each object traces (a) exactly one connected path through space and time, (b) no connected path through space and time, or (c) a path through space and time intersecting another object's path. (Adapted from Spelke and Van de Walle 1993)

Spelke herself has some very definite views about what this understanding consists in. According to Spelke, even very young infants have a theoretical understanding of physical objects and how they behave. Infants are able to represent principles such as those that we have been discussing – the principles of continuity, solidity, and so on. They can use these

principles to make predictions about how objects will behave. They show surprise when those predictions are not met – and lose interest (as measured by looking times) when they are met.

## How Should the Dishabituation Experiments Be Interpreted?

What the infants are doing, according to Spelke (and many others), is not fundamentally different in kind from what scientists do. The infants are making inferences about things that they cannot see on the basis of effects that they can see – just as scientists make inferences about, say, subatomic particles on the basis of trails in a cloud chamber. Infants are little scientists, and the perceptual discriminations that they make reflect their abilities to make inferences about the likely behavior of physical objects; inferences that in turn are grounded in a stored and quasi-theoretical body of knowledge about the physical world – what is sometimes called infant folk physics.

So, what sort of information processing underlies infant folk physics? The physical symbol system hypothesis gives us a natural way of thinking about how rules might be explicitly represented and applied. The idea here is that the basic principles of infant folk physics (such as the principle of continuity) are symbolically represented. These symbolically represented principles allow the infants to compute the probable behavior and trajectory of the objects in the dishabituation experiments. They show surprise when objects do not behave according to the results of the computations.

This view is perfectly consistent with the idea that infant folk physics is importantly different from adult folk physics. Infant folk physics has some puzzling features. Developmental psychologists have found, for example, that infants tend to place more weight on spatiotemporal continuity than on featural continuity. For infants, movement information dominates information about features and properties. Their principal criterion for whether or not an object persists over time is that it should maintain a single trajectory, even if its perceptible properties completely change.

This is why, for example, infants who otherwise perceive differences between the color and form of objects still tend not to show surprise when one object disappears behind a screen and another completely different object emerges at the other side of the screen. For adults, on the other hand, featural constancy is often more important. This is elegantly expressed by the developmental psychologists Alison Gopnik and Andrew Meltzoff:

As adults we individuate and reidentify objects by using both place and trajectory information and static-property information. We also use property information to predict and explain appearances and disappearances. If the same large, distinctive white rabbit appears in the box and later in the hat, I assume it's the same rabbit, even if I don't immediately see a path of movement for it. In fact, I infer an often quite complex invisible path for the object. If I see the green scarf turn into a bunch of flowers as it passes through the conjuror's hand while maintaining its trajectory, I assume it is a



different object. On the other hand, if an object changes its trajectory, even in a very complex way, while maintaining its properties, I will assume it is still the same object.

(Gopnik and Meltzoff 1997: 86)

So, there are some important differences between infant folk physics and adult folk physics. The important point, though, is that for Spelke (and indeed for Gopnik and Meltzoff) both should be understood as theories. Here is how Spelke described her findings in an influential early paper.

I suggest that the infant's mechanism for apprehending objects is a mechanism of thought: an initial *theory* of the physical world whose four principles jointly define an initial *object concept*.

(Spelke 1988: 181)

It is no easy matter to say what a theory actually is, but as Spelke states, the simplest way of thinking about theories is in terms of laws or principles. Laws and principles can be linguistically expressed. This means that they can easily be represented by physical symbol structures. In this respect, thinking about naïve physics as a theory is rather like thinking of grammatical knowledge in terms of rules. In both cases we have cognitive capacities (knowledge of a theory in the one case, and the ability to apply rules in the other) that lend themselves to being modeled in computational terms – as suggested by the physical symbol system hypothesis.

As we saw in the case of grammatical knowledge, however, there are alternatives to this type of computational approach. We can think about knowledge in nonsymbolic ways, exploiting neural network models. Some of the possibilities are sketched out in the next section.

## 11.2

## Neural Network Models of Children's Physical Reasoning

Connectionist modelers have explored alternatives to the theoretical model of infant cognitive abilities that we have just looked at. They have tried to show how a neural network can simulate the behavior of human infants in experiments using the dishabituation paradigm without any relevant principles or rules being explicitly coded into it.

One prominent researcher in this area is the psychologist Yuko Munakata, working with a number of collaborators, including the distinguished connectionist modeler Jay McClelland (who, together with David Rumelhart, edited the two-volume *Parallel Distributed Processing*, which gave such a huge impetus to connectionist approaches to cognitive science). Here is how Munakata and her co-authors describe the basic idea behind their approach, and how it differs from the theoretical model:

Because infants seem to behave in accordance with principles at times, there might be some use to describing their behavior in these terms. The danger, we believe, comes in the tendency to accept these descriptions of behavior as mental entities that are explicitly accessed and used in the production of behavior. That is, one could say that infants'

behavior in a looking-time task accords with a principle of object permanence, in the same way one could say that the motions of the planets accord with Kepler's laws. However, it is a further – and we argue unfounded – step to then conclude that infants actually access and reason with an explicit representation of the principle itself.

The connectionist modelers accept that the dishabituation experiments show that human infants are sensitive to (and react in accordance with) certain basic physical principles (such as the principles of solidity and continuity). But they reject the way that computational theorists interpret this basic fact.

The computational approach and the theoretical model of infant cognition both assume that a cognitive system (whether a human infant, or a computational model) can only act in accordance with, say, the principle of continuity if that principle is explicitly represented in it in a symbolic form. But, according to Munakata and her collaborators, this assumption is wrong – and it can be shown to be wrong by constructing a neural network model that acts in accordance with the principle of continuity even though it does not have that principle symbolically encoded in it. They continue:

We present an alternative approach that focuses on the adaptive mechanisms that may give rise to behavior and on the processes that may underlie change in these mechanisms. We show that one might characterize these mechanisms as behaving in accordance with particular principles (under certain conditions); however, such characterizations would serve more as a shorthand description of the mechanism's behavior, not as a claim that the mechanisms explicitly consult and reason with these principles.

(Munakata et al. 1997: 687)

Their alternative proposal is that infants' understanding of object permanence is essentially practical. The fact that infants successfully perform object permanence tasks does indeed show that they know, for example, that objects continue to exist even when they are not being directly perceived. But this knowledge is not explicitly stored in the form of theoretical principles. In fact, it is not explicitly stored at all. Rather, it is implicitly stored in graded patterns of neural connections that evolve as a function of experience.

The basic idea is that infants' expectations about how objects will behave are driven by patterns of neural activation. These patterns vary in strength due to

- the number of neurons firing
- the strength and number of the connections between them
- the relations between their individual firing rates

So, the types of perceptual sensitivity that we see in the dishabituation paradigms are produced by associative mechanisms of pattern recognition. This is exactly what connectionist networks model so well.

Here is how the process works. When infants observe the "reappearance" of occluded objects, this strengthens the connection between two groups of neurons – between the group of neurons that fire when the object first appears, on the one hand, and the group

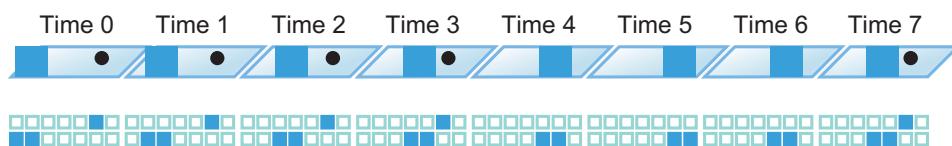


that fires when it reappears, on the other. As a result, the representations of perceived objects (i.e., the patterns of neural activation that accompany the visual perception of an object) persist longer when the object is occluded. So, according to Munakata et al., the infant's "knowledge" of object permanence should be understood in terms of the persistence of object representations, rather than in terms of any explicitly coded principles. This "implicit" understanding of object permanence is the foundation for the theoretical understanding that emerges at a much later stage in development.

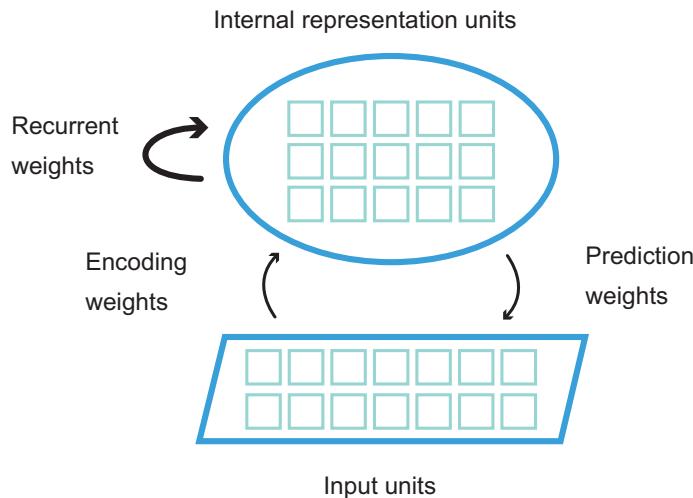
One advantage of their approach is that it explains well-documented behavioral dissociations in infant development. There is good evidence that infants' abilities to act on occluded objects lag a long way behind their perceptual sensitivity to object permanence, as measured in preferential looking tasks. Although perceptual sensitivity to object permanence emerges at around 4 months, infants succeed in searching for hidden objects only at around 8 months. Munakata et al. argue (and their simulations illustrate) that it is possible for a visual object representation to be sufficiently strong to generate expectations about the reappearance of an occluded object, while still being too weak to drive searching behavior.

## Modeling Object Permanence

One of the networks studied by Munakata et al. is designed to simulate a simple object permanence task involving a barrier moving in front of a ball and occluding the ball for a number of time steps. Figure 11.5 shows the inputs to the network as the barrier moves in front of the ball and then back to its original location. The input units are in two rows. The two rows jointly represent the network's "field of view." The bottom layer represents the network's view of the barrier, while the top layer represents the network's view of the ball. As we see in the figure, when the barrier moves in front of the ball there is no input in the ball layer. When the barrier moves to one side, revealing the previously occluded ball, the ball layer is correspondingly activated again. What the network has to do is to learn to represent the ball even when there is no activation in the input layer corresponding to the



**Figure 11.5** A series of inputs to the network as a barrier moves in front of a ball and then back to its original location. The top row shows a schematic drawing of an event in the network's visual field; the bottom row indicates the corresponding pattern of activation presented to the network's input units, with each square representing one unit. Learning in the network is driven by discrepancies between the predictions that the network makes at each time step and the input it receives at the next time step. The correct prediction at one time step corresponds to the input that arrives at the next time step. (Adapted from Munakata et al. 1997)



**Figure 11.6** Recurrent network for learning to anticipate the future position of objects. The pattern of activation on the internal representation units is determined by the current input and by the previous state of the representation units by means of the encoding weights and the recurrent weights, respectively. The network sends a prediction back to the input units to predict the next state of the input. The stimulus input determines the pattern of activation on the input units, but the difference between the pattern predicted and the stimulus input is the signal that drives learning. (Adapted from Munakata et al. 1997)

ball – it needs to find a way of representing the ball even when the ball cannot directly be seen.

In order to design a network that can do this Munakata and her collaborators used a type of network that we have not yet looked at, called a *recurrent network*. These networks are rather different from the feedforward and competitive networks that we have been considering up to now. Like feedforward and competitive networks, they have hidden units whose weights are modified by algorithmic learning rules. But what distinguishes them is that they have a feedback loop that transmits activation from the hidden units back to themselves. This transmission works before the learning rule is applied. This feedback loop allows the network to preserve a “memory” of the pattern of activation in the hidden units at the previous stage.

Figure 11.6 is a schematic representation of their recurrent network. The network has two distinctive features. The first is the set of recurrent weights from the hidden layer back to itself. These function as just described – to give the network information about what happened at the previous temporal stage. The second is a set of connections, with corresponding weights, running from the hidden units to the input units. These weighted connections allow the network to send a prediction to the input units as to what the next set of inputs will be. The network’s learning (which works via the standard backpropagation rule) is driven by the discrepancy between the actual input and the predicted input.

We can think about the network’s “understanding” of object permanence in terms of its sensitivity to the ball’s reappearance from behind the occluder. This sensitivity can in turn



be measured in terms of the accuracy of the network's "prediction" when the ball does eventually reappear. (An accurate prediction is one where the predicted pattern exactly matches the input pattern.) As training progresses the network becomes increasingly proficient at predicting the reappearance of occluded objects over longer and longer periods of occlusion.

What makes this possible is the recurrent connection from the hidden layer back to itself. The activation associated with the "sight" of the ball at a given temporal stage is transmitted to the next stage, even when the ball is not in view. So, for example, at temporal stages 4, 5, and 6 in Figure 11.5, there is no activation in the input units representing the ball. But, once the network's training has progressed far enough, the weights will work in such a way that the memory from the earlier stages is strong enough that the network will correctly predict the reappearance of the ball at temporal stage 7.

How exactly does this work? The researchers found that improved sensitivity to object permanence is directly correlated with the hidden units representing the ball showing similar patterns of activation when the ball is visible and when it is occluded. In effect, they claim, the network is learning to maintain a representation of an occluded object. The network's "understanding" of object permanence is to be analyzed in terms of its ability to maintain such representations. And this comes in degrees. As further simulations reported in the same paper show, a network can maintain representations sufficiently strong to drive perceptual "expectations" but too weak to drive motor behavior. Sensitivity to object permanence is, they suggest, a graded phenomenon – a function of strengthened connections allowing maintained activation patterns – rather than a theoretical achievement.

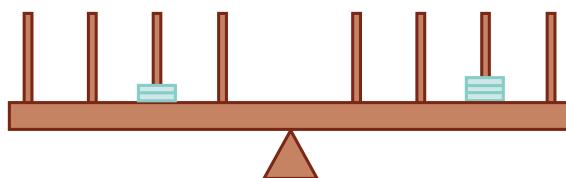


**Exercise 11.2** Explain and assess the significance of this network model for thinking about the information processing underlying object permanence.

## Modeling the Balance Beam Problem

Here is a second example of how connectionist models can provide alternatives to theory-based accounts of infant cognitive development. This is the balance beam problem.

Children are shown a balance beam as in Figure 11.7. The balance beam has a fulcrum and weights at varying distances from the fulcrum. The children are asked whether the beam is in balance and, if not, which side will go down. In different trials the weights are varied, but the children are not given any feedback on whether their answers are correct or not.



**Figure 11.7** A balance beam. Weights can be added at different distances from the fulcrum. Children are asked whether the beam is in balance and, if not, which side will go down.

Research by the developmental psychologist Bob Siegler has shown that children typically go through a series of stages in tackling the balance beam problem – rather like young children learning the past tense of English verbs. And, as in the past tense case, these stages can be summarized in terms of some relatively simple rules. There are four stages and corresponding rules. Siegler identifies these as follows:

*Stage 1* Children think that the side with the greatest number of weights will go down, irrespective of how those weights are arranged. If there are equal numbers of weights on both sides, then the beam is judged to be in balance.

*Stage 2* Now they think that when the weights on each side of the fulcrum are equal, the side on which the weights are furthest away will go down. If this doesn't hold then children either use the first rule or guess.

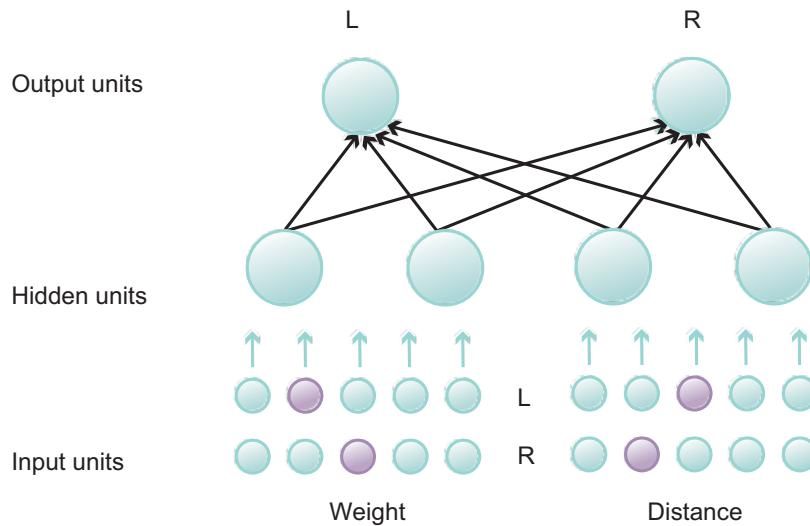
*Stage 3* Now children are able to use the correct rule, understanding that downward force is a function both of weight and of the distance from the fulcrum. But they only manage to do this when the two sides differ in respect *either* to weight *or* to distance, but not both.

*Stage 4* It is usually not until adolescence that children acquire a general competence for balance beam problems – and even then, not all of them do.

The situation here is very similar to the past tense case. And, as in that case, it seems initially plausible to model the child's learning process as a matter of learning a series of rules. But as we've already seen, there are other ways to think about this type of developmental progression. Artificial neural networks offer an alternative way of looking at the phenomenon, illustrating how the appearance of rule-based learning can emerge from a system that does not exploit any explicit rules.

Jay McClelland and E. Jenkins designed an artificial neural network to model children's performance on the balance beam problem. The network is designed to reflect the different types of potential input in solving balance beam-type tasks. The network is illustrated in Figure 11.8. It has four different groups of input units, receiving input about weights and distances for each side of the fulcrum. It is important to realize that the information the network gets is actually quite impoverished. One group of input units will get information corresponding to, say, the weights to be found on one side of the beam. Another group of units will get information corresponding to the distances of those weights from the fulcrum. But these are separate pieces of information. The network needs to work out during training that the two groups of units are carrying information about the same side of the balance beam. The network weights are initially set at random.

As we see in Figure 11.8, the weight units are connected to a pair of hidden units. Likewise, for the distance units. There are no connections between the two pairs of hidden units, but each hidden unit projects to both the output units. The network predicts that the balance beam will come down on the left-hand side when the activation on the left output unit exceeds the activation on the right output unit.



**Figure 11.8** The architecture of the McClelland and Jenkins network for the balance beam problem. (Adapted from Elman et al. 1996)

The McClelland–Jenkins network learns by backpropagation. The discrepancy between the correct output and the actual output on given iterations of the task is propagated backward through the network to adjust the weights of the connections to and from the hidden units.

As the training went on, the network went through a sequence of stages very similar to those that Siegler identified in children. The initial training examples showed much more variation in weight than in distance. This was intended to reflect the fact that children are more used to using weight than distance in determining quantities like overall heaviness. As an artifact of the training schedule, therefore, the network's early discriminations all fell into Siegler's stage 1. As training progressed, the network learned to use distance to solve problems with equal numbers of weights on each side – as per Siegler's stage 2. The final stages of the training saw the network move to Siegler's stage 3, correctly using both weight and distance provided that the two sides differed only on one dimension, but not on both. The McClelland–Jenkins network did not arrive at Siegler's stage 4. But a similar network designed by Jay McClelland did end up showing all four stages.

The moral to be drawn from this example is rather similar to the moral of the tense learning networks we looked at in Section 10.2. Like tense learning, progress on the balance beam problem can be characterized as a step-like progression. Each step seems to involve exploiting a different rule. The most natural way of modeling this kind of learning pattern would be via a model that had these rules explicitly wired into it – exactly the sort of model that would be suggested by the physical symbol system hypothesis. The qualitative progression between different stages would be explained by the transition from one rule to another.

Neural network models show us, however, that step-like progressions can emerge without explicit rules. There are only two rules explicitly programmed into the network – the activation rule governing the spread of activation forward throughout the network, and the backpropagation rules governing the spread of error backward through the network. There is nothing in the network corresponding to the rules in terms of which it might be described. Nor are there any sharp boundaries between the type of learning at different stages, even though its actual performance on the task has a clearly identifiable step-like structure.

## 11.3 Conclusion: The Question of Levels

The models we've been looking at have revealed some of the great strengths of artificial neural networks – particularly when it comes to modeling complicated learning trajectories. We have seen how representations in neural networks are distributed across different hidden units, and how hard it can be to find any sort of straightforward mapping between what is going on inside the network and the task that the network is performing. In this final section we will step back from the details of individual neural network models to look briefly at a very important concern that some cognitive scientists have raised about the whole enterprise of neural network modeling.

To set the scene, think back to David Marr's tri-level hypothesis, which we looked at in Section 2.3. Marr distinguished three different levels at which cognitive scientists can think about a given cognitive system. A quick reminder:

- *The computational level* provides a general characterization of the information-processing task that the system is trying to perform.
- *The algorithmic level* identifies a particular algorithm or set of algorithms that can carry out the task identified at the computational level.
- *The implementational level* explains how the algorithm is actually realized in the system.

Think about a Turing machine, for example. Analysis at the computational level identifies the general information-processing task that the machine is performing – e.g., computing the arithmetical function of addition. An analysis at the algorithmic level will come up with a specific machine table that will compute this function. And then, when we turn to the implementational level, what we are interested in is how to build a physical system that will run that algorithm.



**Exercise 11.3** Pick an example of a cognitive phenomenon from anywhere in this book and use it to illustrate in your own words the difference between Marr's three levels of analysis.

The difference between the algorithmic and implementational levels is very important. The implementational level is the level of engineering and machinery. In contrast, the algorithmic level is the level of discrete information-processing steps, each governed by specific rules. Our Turing machine might take the form of a digital computer. In this case the algorithmic-level



analysis would correspond to the program that the computer is running, while the implementational analysis would explain how that program is realized in the hardware of the computer.

Physical symbol theorists have tended to be very explicit about the level at which their accounts are pitched. As one would expect, given the emphasis on algorithms and rules for manipulating symbol structures, the physical symbol system hypothesis is aimed squarely as an algorithmic-level account. It is not an engineering-level account of information-processing machinery. Rather, it needs to be supplemented by such an account.

So, how should we think about artificial neural networks? If we have a connectionist model of, say, past tense learning, should we think about it as an algorithmic-level account? Or should we think about it as an account offered at the implementational level? Do artificial neural networks tell us about the abstract nature of the information-processing algorithms that can solve particular types of cognitive task? Or do they simply give us insight into the machinery that might run those information-processing algorithms?

This is important because artificial neural networks will only count as alternatives to physical symbol systems if they turn out to be algorithmic-level accounts. The whole contrast that we have been exploring in the last two chapters between neural network models of information processing and physical symbol system models depends upon understanding neural networks at the algorithmic level.

A number of physical symbol theorists (most prominently Jerry Fodor and Zenon W. Pylyshyn) have used this point to make a powerful objection to the whole enterprise of artificial neural network modeling. In effect, their argument is this. We can think about artificial neural networks either at the implementational or at the algorithmic level. If we think about them at the implementational level, then they are not really an alternative to the physical symbol system hypothesis at all. They are simply offering models of how physical symbol systems can be implemented.

But, Fodor and Pylyshyn argue, artificial neural networks shouldn't be seen as algorithmic-level accounts. As language of thought theorists (as described in Chapter 4), they think that cognition can only be understood in terms of the rule-governed transformation of abstract symbol structures – a manipulation that is sensitive only to the formal, syntactic features of those symbol structures. And this only works when we have symbol structures composed of separable and recombinable components.

But artificial neural networks simply do not have separable and recombinable components. They have a very different kind of structure. The state of a network at any given moment is fixed by the state of all the units it contains. And since each distinct unit has a range of possible activation levels, there are as many different possible dimensions of variation for the network as a whole as there are units. So, if there are  $n$  such units, then we can think of the state of the network at any given moment as a point in an  $n$ -dimensional space – standardly called the *activation space* of the system.

Computation in an artificial neural network is really a movement from one point in the network's activation space to another. But when you think about it like that, it is hard to see how the notion of structure can apply at all. A point on a line does not have any

structure. It does not have separable and recombinable components. Nor does a point on the plane (i.e., in two-dimensional space) – or a point in any  $n$ -dimensional space.

This allows us to see the force of Fodor and Pylyshyn's argument. We can put it in the form of a dilemma. Either neural networks contain representations with separable and recombinable components, or they do not. If they do contain such representations, then they are not really alternatives to the physical symbol system hypothesis. In fact, they will just turn out to be ingenious ways of implementing physical symbol systems. But if, on the other hand, they do not contain such representations, then (according to Fodor and Pylyshyn) they have absolutely no plausibility as algorithmic-level models of information processing. Here is the argument, represented schematically:

- 1 Either artificial neural networks contain representations with separable and recombinable components, or they do not.
- 2 If they do contain such representations, then they are simply implementations of physical symbol systems.
- 3 If they do not contain such representations, then they cannot plausibly be described as algorithmic information processors.
- 4 Either way, therefore, artificial neural networks are not serious competitors to the physical symbol system hypothesis.

This argument is certainly elegant. You may well feel, though, that it is begging the question. After all, the whole point of the neural network models we have been looking at in this chapter (and Chapter 10) has been to try to show that there can be information processing that does *not* require the type of rule-governed symbol manipulation at the heart of the physical symbol system hypothesis. In a sense, the models themselves are the best advertisement for artificial neural networks as genuine alternative models of information processing – rather than simply implementations of physical symbol systems.



#### **Exercise 11.4 Assess in your own words step 3 in Fodor and Pylyshyn's argument.**

In any case, there is no law that says that there is only one type of information processing. Perhaps the physical symbol system approach and the neural networks approach can co-exist. It may turn out that they are each suitable for different information-processing tasks. When we explored the language of thought hypothesis, for example, we placed considerable emphasis on the role of propositional attitudes such as belief and desire in causing behavior. The interplay of syntax and semantics in the language of thought was intended to capture the idea that beliefs and desires could bring about behavior in virtue of how they represent the world. But the types of task we have been looking at in these last two chapters seem very different – closer to perception and pattern recognition than to abstract symbol manipulation. It may turn out that different types of cognitive task require fundamentally different types of information processing.



## Summary

This chapter first reviewed experiments indicating infants' expectations of physical objects and their behavior, resulting in a theory-like folk physics being attributed to infants. We then looked at a neural network model of infant understanding of object permanence. What we see is that, without really having any rules encoded in it, the network can accurately model infant expectations about when occluded objects are going to reappear from behind a barrier. It learns to do this because of feedback connections from the hidden units to themselves, which function as a type of memory within the network. The chapter ended by considering a famous dilemma that Fodor and Pylyshyn have posed for neural network models.

## Checklist

### Object Permanence in Infancy

- (1) According to the traditional view, the perceptual universe of the infant is a "blooming, buzzing, confusion" with infants only coming to understand object permanence (i.e., that objects continue to exist when they are not directly perceived) at the age of 8 months or so.
- (2) Recent studies using the dishabituation paradigm have led many developmental psychologists to revise this view and to claim that even very young infants inhabit a highly structured and orderly perceptual universe.
- (3) Researchers such as Elizabeth Spelke have argued that young infants are able to parse the visual array into objects that behave according to certain basic physical principles. This is often called object permanence.

### Neural Network Models of Object Permanence

- (1) Many cognitive scientists think that object permanence depends upon computations that exploit explicitly represented physical principles (a primitive folk physics)
- (2) Munakata's neural network model suggests that object permanence might be a matter of having representations of objects that persist when the object is occluded, rather than explicitly representing physical principles.
- (3) This network is a recurrent network, where a feedback loop from the hidden units back to themselves functions as a type of memory.

### Modeling the Balance Beam Problem

- (1) Studies have explored how young children reason about what will happen to objects of different weights placed on opposite sides of a balance beam at different distances from the fulcrum.
- (2) Children typically go through a fairly standard learning trajectory, where weight and distance are differently weighted at different stages, before they understand that downward force is a function of both weight and distance
- (3) This trajectory can be modeled by a neural network whose hidden units develop associative connections between inputs corresponding to the weights and distances.

- (4) This suggests that we can think about children's developing physical understanding without assuming that they are developing an increasingly sophisticated theoretical understanding of physical principles connecting force, weight, and distance.

### The Fodor–Pylyshyn Objection to Neural Network Modeling

- (1) Fodor and Pylyshyn start with a dilemma: Either artificial neural networks contain representations with separable and recombinable components, or they do not.
- (2) If neural networks do contain such representations, then (they argue) the networks are simply implementations of physical symbol systems.
- (3) But if they do not contain such representations, then (according to Fodor and Pylyshyn) they cannot plausibly be described as algorithmic information processors.
- (4) Either way, Fodor and Pylyshyn argue, artificial neural networks are not serious competitors to the physical symbol system hypothesis.
- (5) *But* – this seems to be begging the question, since the central claim of the neural networks is that information processing need not require the type of rule-governed symbol manipulation at the heart of the physical symbol system hypothesis.

## Further Reading

The drawbridge experiments described in Section 9.3 were first present in Baillargeon 1986 and 1987. They have been extensively discussed and developed since then. For a more recent model, see Wang and Baillargeon 2008. Spelke's experiments using the dishabituation paradigm are reviewed in many places (e.g., Spelke et al. 1995). A general discussion of habituation methodology can be found in Oakes 2010. Spelke and Kinzler 2007 reviews evidence for infant "core knowledge" in understanding objects, actions, number, and space. Susan Carey and Renée Baillargeon have extended Spelke's "core knowledge" in a number of ways. Summaries can be found in Carey and Spelke 1996, Carey 2009, Baillargeon et al. 2010, and Baillargeon and Carey 2012. Woodward and Needham 2009 is a collection of review articles on infant cognition. Hespos and van Marle 2012 provide a summary pertaining specifically to infants' knowledge of objects. Cacchione 2013 continues experimental work on infant perception of cohesion. The "child as little scientist" theory is engagingly presented in Gopnik and Meltzoff 1997.

One of the first papers exploring connectionist approaches to object permanence was Mareschal, Plunkett, and Harris 1995. See further Mareschal and Johnson 2002. The papers discussed in the text are Munakata et al. 1997, Munakata 2001, and Munakata and McClelland 2003. For a book-length treatment of the power of connectionist approaches in thinking about cognitive development, see Elman et al. 1996 – which also contains a detailed account of the balance beam network discussed in Section 9.4 (originally presented in McClelland and Jenkins 1991). Plunkett and Elman 1997 is an accompanying workbook with software. Marcus 2003 attempts to integrate connectionist and symbolic approaches. Elman 2005 is another good review. A critical view can be found in Quinlan et al. 2007.



The Fodor and Pylyshyn argument discussed in Section 9.5 can be found in Fodor and Pylyshyn 1988. It has been widely discussed. A number of important papers are collected in Macdonald and Macdonald 1995. See chapter 9 of Bermúdez 2005 for a general discussion and further references and Calvo and Symons 2014 for a book-length overview. Jansen and Watter 2012 proposes a network that the authors claim displays strong systematicity.





## CHAPTER TWELVE

# Machine Learning: From Expert Systems to Deep Learning

### OVERVIEW 307

#### 12.1 Expert Systems and Machine Learning 308

Expert Systems and Decision Trees 308  
ID3: An Algorithm for Machine Learning 310

#### 12.2 Representation Learning and Deep Learning 315

Deep Learning and the Visual Cortex 318

#### 12.3 The Machinery of Deep Learning 321

Autoencoders 322  
Convolutional Neural Networks 324  
Sparse Connectivity 325  
Shared Weights 326  
Invariance under Translation 326

#### 12.4 Deep Reinforcement Learning 327



## Overview

This chapter is dedicated to machine learning, one of the hottest topics in contemporary AI and the key to the success of multi-billion-dollar corporations such as Google, Facebook, and Amazon.

We begin in Section 12.1 by introducing the idea of expert systems, computer programs that are designed to replicate (and improve on) the performance of human experts in specialized domains, such as identifying diseases in humans and plants, or processing credit card applications. These programs can often be represented as decision trees. There are different ways of constructing expert systems, however. One way is to start with human experts and write a program that codifies their collective knowledge. Alternatively, machine learning algorithms can be used to construct a decision tree by analyzing large databases of examples and deriving rules that can then be used to classify new examples. We illustrate this through ID3, which is an example of a traditional machine learning algorithm.

Traditional algorithms such as ID3 are still highly dependent upon how their databases are labeled and constructed. They typically require lengthy and complex processes of *feature*

*engineering* to structure the databases. In the subfield of machine learning known as *representation learning* computer scientists write programs that will do their own feature engineering. And within representation learning as a whole, the greatest advances have come from what is called *deep learning*.

We explore deep learning in Sections 12.2, 12.3, and 12.4. In Section 12.2 we see how deep learning programs typically involve multiple layers of artificial neural networks. These networks are hierarchically organized to extract increasingly complex information from the raw data. The mammalian visual system is an explicit inspiration for this type of construction. Section 12.3 introduces two representative types of neural networks used for deep learning – autoencoders and convolutional neural networks. Finally, in Section 12.4 we look at the reinforcement learning methods behind two of the most spectacular examples of deep learning – the AlphaGo and AlphaGo Zero programs created by Google’s Deep Mind research team.

## 12.1

### Expert Systems and Machine Learning

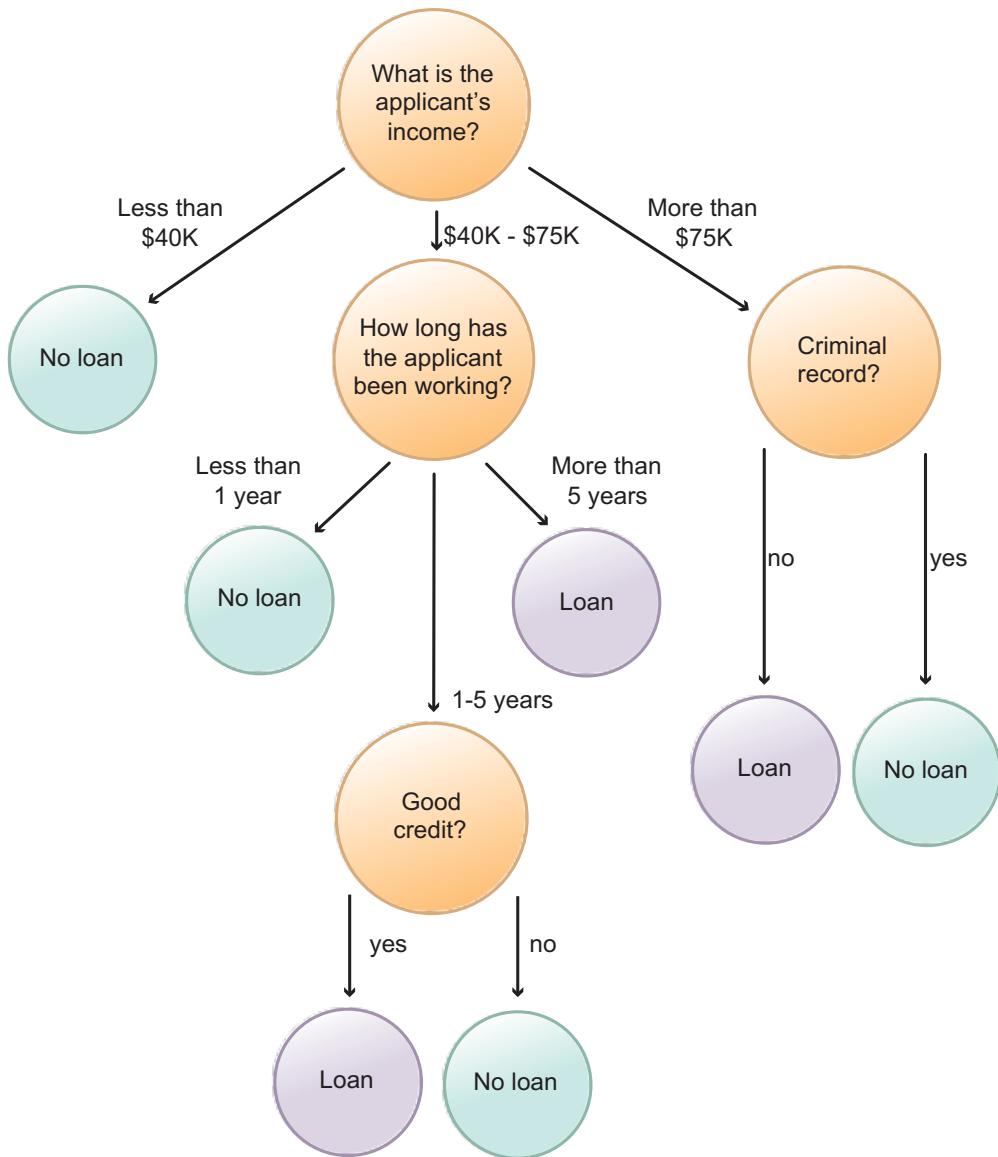
In the field of AI known as *expert systems* research, researchers write computer programs to reproduce (and ideally improve on) the performance of human beings who are expert in a particular domain.

Expert systems programs are typically applied in narrowly defined domains to solve very determinate problems, such as diagnosing specific medical disorders. A well-known expert systems program called MYCIN was developed at Stanford University in the early 1970s. MYCIN was designed to simulate a human expert in diagnosing infectious diseases. It took in information from doctors on a particular patient’s symptoms, medical history, and blood tests, asking for any required information that it did not already have. It then analyzed this information using a knowledge base of about 600 heuristic rules about infectious diseases derived from clinical experts and textbooks.

MYCIN produced a number of different diagnoses and recommendations for antibiotic treatments. It was able to calculate its degree of confidence in each diagnosis and so present its findings as a prioritized list. Although MYCIN was never actually used as the sole tool for diagnosing patients, a widely reported study at Stanford University’s medical school found that it produced an acceptable diagnosis in 69 percent of cases. You may think that 69 percent is not very high, but it turns out to be significantly higher than infectious disease experts who were using the same rules and information.

### Expert Systems and Decision Trees

Expert systems have become very deeply entrenched in the financial services industry, particularly for mortgage loan applications and tax advice. Most banks these days have online “wizards” that will take mortgage applicants through a series of simple questions designed to lead to a decision on the applicant’s “mortgage-worthiness.” Mortgage wizards can be represented through *decision trees*. In the simplest form of decision tree each node



**Figure 12.1** A decision tree illustrating a mortgage expert system. (From Friedenberg and Silverman 2006)

corresponds to a question. Each node has several branches leading from it, corresponding to different answers to the question.

Figure 12.1 illustrates a very simple schematic expert system for a loan decision tree. Two features of this decision tree are worth highlighting. First, it offers a fixed decision procedure. Whatever answers the loan applicant gives to the fixed questions, the decision tree will eventually come up with a recommendation. Second, the presentation in tree form is completely inessential. We can easily convey what is going on in terms of explicit rules, such as the following:

IF income less than \$40K THEN no loan

IF income greater than \$75K AND no criminal record THEN loan

IF income between \$40K and \$75K AND applicant working for 1–5 years AND credit not good THEN no loan

(I have used uppercase letters to bring out the logical structure of the rules.) When the decision tree is written as a computer program it may well be written using explicit rules such as these.

The decision tree works because of the questions that are asked at each node. When taken together the questions exhaust the space of possibilities. Each question partitions the possibility space in such a way that each branch of the tree leads to a unique outcome (what computer scientists call a *terminal leaf or node*). But how are we supposed to get to these questions? How does the decision tree get designed, as it were?

One possibility would be to ask a team of mortgage loan officers to sit down and work out a decision tree that captures the practices at their bank. This could then be used as the basis for writing a program in a suitable programming language. This would be fine, and it is no doubt how many expert systems programs are actually written (particularly in the mortgage area). But from the perspective of AI this would not be very interesting. It would be an expert system only in a very derivative sense. The real expert system would be the team of mortgage loan professionals. Much more interesting would be a program that was capable of producing its own decision tree – a program capable of imposing its own structure upon the problem and working out what would count as a solution. How would this work?

Suppose that we have a huge database of all the loan decisions that the bank has taken over a long period of time, together with all the relevant information about the applicants – their income, work history, credit rating, and so on. If we can find a way of representing the bank’s past decisions in the form of a decision tree, so that each branch of the tree ends either in the loan being given or the loan being declined, then we can use that decision tree to “process” new applications.

This is a classic example of the type of problem tackled in the branch of AI known as *machine learning* (a subfield in expert systems research). The challenge is to produce an algorithm that will organize a complex database in terms of some attribute we are particularly interested in (such as an applicant’s loan-worthiness, in the example we are considering). The organization takes the form of a decision tree, which will determine whether or not the attribute holds in a given case (i.e., whether or not the applicant is loan-worthy).

## ID3: An Algorithm for Machine Learning

This section explores an influential machine learning algorithm developed by the computer scientist Ross Quinlan. Quinlan developed the ID3 learning algorithm while working at the University of Sydney in Australia. He now runs a company called RuleQuest Research which is commercially marketing updated and more efficient versions of the ID3 algorithm.

A machine learning algorithm works on a vast database of information. It looks for regularities in the database that will allow it to construct a decision tree. Machine learning algorithms such as ID3 only work on databases that take a very specific form.



The basic objects in the database are standardly called *examples*. In the loan application decision tree that we looked at earlier, the examples are loan applicants. These loan applicants can be classified in terms of a certain number of *attributes*. Each example has a value for each attribute. So, for example, if the attribute is *Credit History?*, then the possible values are *Good or Bad* and each mortgage applicant is assigned exactly one of these values. The attribute we are interested in is the *target attribute*. In our example the target attribute is *Loan* and the two possible values are *Yes and No*. Again, every applicant is either offered a loan or is turned down.

The attributes work to divide the examples into two or more classes. So, for example, the attribute at the top of the decision tree is *Income?*. This attribute divides the loan applicants into three groups. As we move down each branch of the tree each node is an attribute that divides the branch into two or more further branches. Each branch ends when it arrives at a value for the target attribute (i.e., when the decision is made on whether to give the loan or not).

The ID3 algorithm exploits the basic fact that each attribute divides the set of examples into two or more classes. It assigns attributes to nodes, identifying, for each node in the decision tree, which attribute would be most informative at that point. That is, it identifies at each node which attribute would divide the remaining examples up in the most informative way.

The ID3 algorithm uses a statistical measure of informativeness, standardly called *information gain*. Information gain measures how well a particular attribute classifies a set of examples. At each node the algorithm chooses the remaining attribute with the highest information gain.

The concept of information gain is itself defined in terms of a more fundamental measure called *entropy*. (Warning: You may have come across the concept of entropy in physics, where it features for example in the second law of thermodynamics. Entropy is defined somewhat differently in information theory than in physics and it is the information-theoretic use that we are interested in here.) We can think of entropy as a measure of uncertainty.

Once we have a formula for calculating entropy we can calculate information gain relative to a particular attribute, and then we can use that to construct the decision tree. Basically, for each attribute, the algorithm works out how well the attribute organizes the remaining examples. It does this by calculating how much the entropy would be reduced if the set were classified according to that attribute. This gives a measure of the information gain for each attribute. Then the algorithm assigns the attribute with the highest information gain to the first node on the tree. And the process is continued until each branch of the tree ends in a value for the target attribute.

To see how ID3 works, consider this relatively simple problem – deciding whether or not the weather is suitable for playing tennis. Imagine that, as keen tennis players who seriously consider playing tennis every day, we collect information for 2 weeks. For each day we log the principal meteorological data and note whether or not we decide to play tennis on that day. We want to use this information to construct a decision tree that we can use in the future.

So, the target attribute is *Play Tennis?* Here are the other attributes with the values they can take.

*Outlook?* {sunny, overcast, rain}

*Temperature?* {hot, mild, cool}

Humidity?	{high, low, normal}
Wind?	{weak, strong}

And here is our database.

DAY	OUTLOOK?	TEMPERATURE?	HUMIDITY?	WIND?	PLAY TENNIS?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Even this relatively small database is completely overwhelming. It is very hard to find any correlations between the target attribute and the other attributes. Fortunately, though, this is exactly the sort of problem that ID3 can solve.

The first step is to find the attribute with the highest information gain. When ID3 calculates the information gain for all four attributes the results come out as follows:

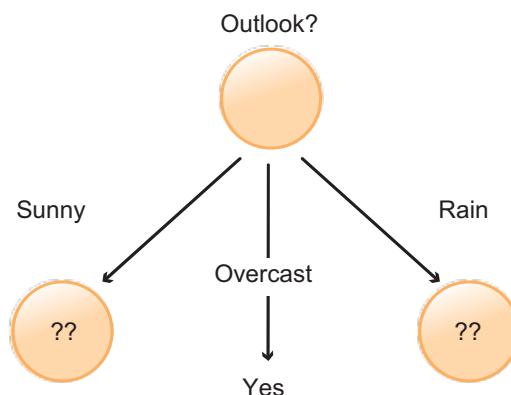
Gain (S, Outlook?)	=	0.246
Gain (S, Temperature?)	=	0.029
Gain (S, Humidity?)	=	0.151
Gain (S, Wind?)	=	0.048



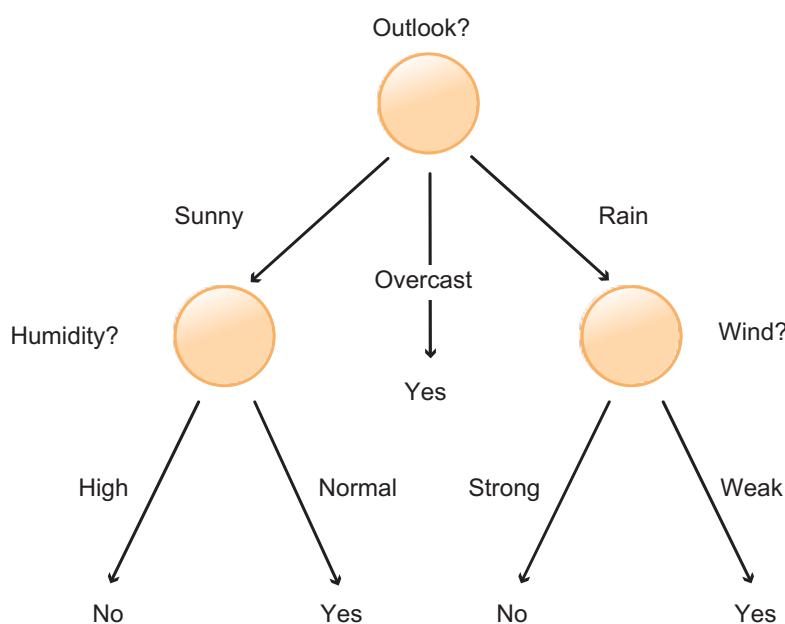
So, it is clear what ID3 will do. The information gain is highest for *Outlook?* and so that is the attribute it assigns to the first node in the decision tree. The decision tree looks like the one in Figure 12.2.

Each of the three branches coming down from the first node corresponds to one of the three possible values for *Outlook?*. Two of the branches (Sunny and Rain) lead to further nodes, while the middle branch immediately ends.

It turns out that assigning attributes to these two nodes is all that is required for a comprehensive decision tree – i.e., for a decision tree that will tell us whether or not to play tennis in any combination of meteorological conditions. The final decision tree is illustrated in Figure 12.3.



**Figure 12.2** The first node on the decision tree for the tennis problem. *Outlook* is the first node on the decision tree because it has the highest information gain.



**Figure 12.3** The complete decision tree generated by the ID3 algorithm.



### Exercise 12.1 Check that this decision tree works for two of the lines in the original database.

This is a “toy” example. But there are plenty of “real-life” examples of how successful ID3 can be. Here is one.

In the late 1970s Ryszard Michalski and Richard Chilausky, two computer scientists at the University of Illinois (deep in the agricultural heartland of America’s Midwest), used ID3 to devise an expert system for diagnosing diseases in soybeans, one of Illinois’s most important crops. This is a rather more difficult problem, since there are nineteen common diseases threatening soybean crops. Each disease is standardly diagnosed in terms of clusters of thirty-five different symptoms. In this case, therefore, the target attribute has nineteen different possible values and there are thirty-five different attributes. Many of these attributes also have multiple possible values.

In order to appreciate how complicated this problem is, look at Figure 12.4. This is part of a questionnaire sent to soybean farmers with diseased crops. It asks for details on a wide range of attributes. Completed questionnaires such as this one were one of the inputs to the initial database. They were supplemented by textbook analyses and lengthy

<b>ENVIRONMENTAL DESCRIPTORS</b> TIME OF OCCURRENCE = ? PLANT STAND = ? PRECIPITATION = ? TEMPERATURE = ? OCCURRENCE OF HAIL = ?
<b>PLANT GLOBAL DESCRIPTORS</b> SEVERITY = ? SEED TREATMENT = ? PLANT HEIGHT = ?
<b>PLANT LOCAL DESCRIPTORS</b> CONDITION OF LEAVES = ? LEAFSPOTS-HALOS = ? LEAFSPOTS-MARGINS = ? LEAFSPOT SIZE LEAF SHREDDING = ? LEAF MALFORMATION = ? LEAF MILDEW GROWTH

**Figure 12.4** A sample completed questionnaire used as input to an ID3-based expert system for diagnosing diseases in soybean crops. (Adapted from Michalski and Chilausky 1980)



consultations with a local plant pathologist. The total database on which ID3 was trained comprised 307 different examples.

Michalski and Chilausky were interested not just in whether ID3 could use the training examples to construct a decision tree. They wanted to compare the resulting decision tree to the performance of a human expert. After all, what better gauge could there be of whether they really had succeeded in constructing an expert system? And so, they tested the program on 376 new cases and compared its diagnoses to those made by various experts on plant disease (including the author of the textbook that they had originally used to compile the database). As it turned out, the expert system did much better than the human expert on the same 376 cases. In fact, it made only two mistakes, giving it a 99.5 percent success rate, compared to the 87 percent success rate of the human experts.

## 12.2

## Representation Learning and Deep Learning

A machine learning algorithm organizes a complex database in terms of one or more target attributes. The objects in the database (the examples) are labeled in terms of a much larger number of attributes. And then, as we just saw with the weather example, the learning algorithm works how to classify the examples progressively in terms of the nontarget attributes until an answer is reached on the target attributes. We can represent this process as a decision tree.

This is very different from a decision tree that is reached by pooling and organizing the insights of a team of experts. Imagine the work that it would take, for example, to get a team of tennis coaches to organize the little database from the last section into a simple set of rules that could tell you whether or not to play tennis in different weather conditions. And, as we saw from the soybean disease example, expert systems trained with machine learning algorithms can outperform human experts.

But still, think about all the work that has to be done before the machine learning algorithm can do its job. When Michalski and Chilausky used ID3 to train an expert system to diagnose diseases in soybean crops, they started with a database that was already highly organized. As shown in Figure 12.4, the examples in the database were classified in terms of highly complex attributes, including different types of information about the plant's leaves, fruit pods, and seeds, as well as general facts about the weather.

In a sense, therefore, much of the work had already been done before the machine learning algorithm had even got going. And so, you might think, all of this preliminary work really reduces the extent to which the algorithm can properly be described as learning. Applying a machine learning algorithm such as ID3 to a database depends crucially upon how the examples in the database are labeled and categorized. And that is an activity that is typically performed by human experts.

There are many people who think that the whole idea of artificial intelligence is confused. So, it is often said, for example, that computers can only follow rules blindly. They are not intelligent themselves. They merely borrow the intelligence of their programmers. John Searle's Russian room argument (see Section 4.3) is a classic expression of this

point of view. And you can see how people unsympathetic to AI for these reasons would think about machine learning algorithms such as ID3. ID3 is impressive in a very limited sense, they might say, but the real intelligence comes from the human experts who designed the questionnaires and wrote the textbooks on which the questionnaires are based.

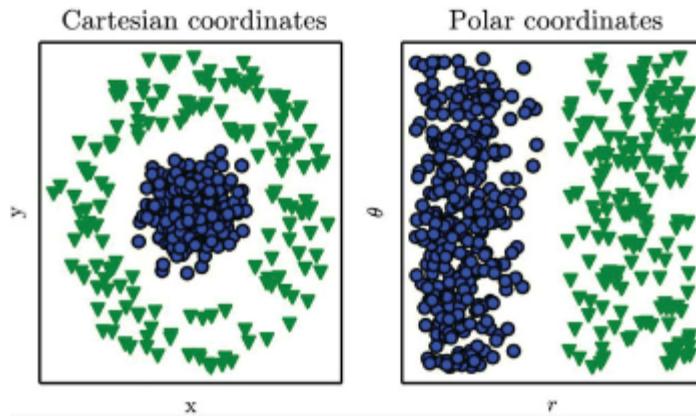
The general thought here is that computers cannot work with raw data such as photographs of diseased soybean plants, for example. They can only work with data that has already been interpreted. Something like this is in fact the guiding idea behind a very familiar experience from browsing the internet. You have probably had plenty of experience with Google's CAPTCHA tool for detecting bots on the internet. The name, CAPTCHA, is an acronym for Completely Automated Public Turing Test to Keep Computers and Humans Apart. A typical image-based CAPTCHA might present you with a photograph divided into 9 subimages and ask you to click on all the subimages that, say, contain part of a car. This is a task than humans can perform with ease, but computers have long had trouble with.

Imagine trying to write a program that will pick out cars in photographs. You know that cars typically have four wheels, but it would be a very unusual photograph that revealed all four of them. And a car propped up on bricks without wheels is still a car. Nor is it easy to define a wheel. We know that wheels are round, but that doesn't help when you can only see part of it. Nor when you can only see it from an angle. The more one thinks about it, the harder it seems to program a computer than can solve CAPTCHA-type tasks – which is why the CAPTCHA widgets are so good at their job.

So, you can think of machine learning algorithms such as ID3 as operating on data that has already been put through a CAPTCHA-like process by human interpreters. The data is labeled with the features that the algorithm will work on. In machine learning this is called *feature engineering*. It is not a trivial task, because data sets can typically be labeled in many different ways and some ways are much easier to work with than others.

For a simple example, consider elementary arithmetic. You probably learned some simple algorithms for multiplying numbers at a very early age. These algorithms are so familiar that it is easy to forget how dependent they are on representing numbers in base 10. You should have no difficulty multiplying 43 by 17. But now try doing the same thing in Roman notation (where you now have to multiply XXXXIII by XVII). Or in binary notation (where the task is now to multiply 101011 by 10001). Computers typically work in binary, but you will probably struggle with binary multiplication.

For a more complicated example, look at the two diagrams in Figure 12.5. Imagine that you have a database of examples that you are trying to divide into two categories. One way to do that is to use the information in the database to plot all of the examples on a scatterplot graph and then see if you can write down the equation of a line that separates the two groups. But you have a choice of coordinate systems. In the left-hand diagram, the examples are all plotted using Cartesian coordinates (i.e., in terms of their position on an ordinary  $x$ -axis and  $y$ -axis). As you can see, the examples do clearly fall into two groups, but since the groups are concentric bands you will need to write the equation for a circle.



**Figure 12.5** Different ways of distinguishing two groups in a database of examples. The left-hand representation uses Cartesian coordinates, while the right-hand representation uses polar coordinates. The right-hand representation makes it much easier, for example, to write the equation for a line separating the two groups.

Exactly the same data is presented in the right-hand diagram, but here they are plotted using polar coordinates (e.g., in terms of their distance from a reference point in the middle of the two concentric circles and their angle from the reference direction of straight upward). Now the two categories of examples are separated by a straight line and the equation is much easier to write.



**Exercise 12.2** Can you think of another example where the difficulty of solving a problem varies according to how the problem is formulated?

These two illustrations are both examples of feature engineering – coding the examples in the database in terms of features that will make it easier to solve the relevant problem. For the multiplication example, the key feature of the numbers being multiplied are how they can be represented in our ordinary base 10 notation (what we more often call decimal notation), as opposed to how they can be represented in base 2 notation (binary) or the Roman version of base 10 notation. And for the scatterplot example, the key features are how the examples can be represented in terms of their distance and angle from a reference point and reference direction.

Standard machine learning algorithms such as ID3 can only get to work once the feature engineering has been carried out, typically by programmers and other human experts. And this feature engineering is often the most complicated and time-consuming part of the process. It is not hard to see the advantages of designing programs and algorithms that would be able to do this feature engineering for us. The technical name for this is *representation learning* or *feature learning*.

This brings us to *deep learning*, which will occupy us for the rest of this chapter. Deep learning is one of the hottest topics in contemporary AI and computer science, pioneered by university researchers and by teams at major tech companies (such as the Deep Mind

team associated with Google). Deep learning algorithms have been responsible for important theoretical advances in areas of AI that had traditionally been viewed as most challenging and intractable. These include computer vision and natural language processing, areas where computer scientists have long struggled to make anything but the most basic advances.

Deep learning has also had practical applications across a wide range of areas. These include software for reading handwriting (not just in languages with alphabets, such as English, but in Chinese) and hence for processing handwritten documents, such as checks, for example. Deep learning has been extensively applied to speech recognition. As a result, voice-based AI, such as Amazon's Alexa, is typically powered by deep learning algorithms. Deep learning is also a key component in the technology behind self-driving cars.

As you might expect, there are many different understandings of deep learning. But here is a very useful characterization from a review article in the journal *Nature* by three pioneers of deep learning – Yann LeCun, Yoshua Bengio, and Geoffrey Hinton:

Representation learning is a set of methods that allows a machine to be fed with raw data and to automatically discover the representations needed for detection or classification. Deep-learning methods are representation-learning methods with multiple levels of representation, obtained by composing simple but non-linear modules that each transform the representation at one level (starting with the raw input) into a representation at a higher, slightly more abstract level.

(LeCun, Bengio, and Hinton 2015)

For LeCun, Bengio, and Hinton, therefore, deep learning is a special type of representation learning. We will start by unpacking this general definition, focusing in particular on the general idea of multiple and hierarchical levels of representation. We will then look at specific architectures and algorithms for deep learning.

## Deep Learning and the Visual Cortex

Deep learning systems are typically constructed from multiple layers of artificial neural networks. In that respect, deep learning is already rather different from more traditional machine learning, which does not usually use artificial neural networks. While it is a descendant of the connectionist models that we have explored in earlier chapters, the design of the individual networks and the overall architecture of the systems is very different from anything that we have looked at up to now.

Deep learning theorists often explicitly appeal to the mammalian visual cortex as a model and inspiration. This is for two reasons. First, the mammalian visual cortex is itself a great example of representation learning. And second, the general design of the mammalian visual cortex is a blueprint for designing representation learning algorithms, particularly because it is hierarchically organized, with different areas interpreting progressively more complex and structured representations.



You can think about the visual system as confronting a classic problem of representation learning. The visual system has to take a complex pattern of unstructured stimuli in the visual field and interpret them into representations that can then serve as input to more complex cognitive functions, such as object recognition. It is a natural representation learning system. By analogy, an artificial representation learning system has the parallel task of taking complex and unstructured raw data and transforming it into representations that can serve as inputs to systems carrying out more complex tasks, such as classification.

The organization and functioning of the mammalian visual cortex is relatively well understood. There is a strong consensus among cognitive neuroscientists that information in the visual cortex is processed hierarchically. Information flows through a progression of different areas, each of which generates representations of increasing complexity.

The first station is the *lateral geniculate nucleus* (LGN), which receives input directly from the retina. In fact, neurons in LGN have similar receptive fields to neurons in the retina. So, really, you can think of LGN as a relay – as the gateway to the visual cortex. LGN projects to area V1 (also known as the primary visual cortex, which maps more or less onto the anatomically defined area known as the striate cortex, because of its conspicuous stripe). Area V1 is where information processing proper begins.

Neurons in V1 are sensitive to low-level features of the visual field, such as orientation and direction of movement. Like LGN, V1 is retinotopically organized – i.e., neighboring regions of the visual field are represented by neighboring regions of V1. In broad outline, V1 takes the retinotopic map coming from LGN and filters it in a way that accentuates edges and the contours. This is the first step in parsing the raw data arriving at the retina into objects.

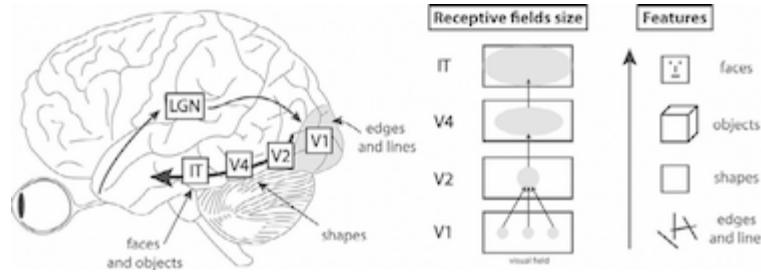
Area V1 projects to area V2 (also known as the secondary visual cortex). Neurons in V2 are tuned to the same features as neurons in V1, as well as to more complex features, such as shape and depth. V2 is much less well understood than V1, but neurons there are sensitive to multiedge features such as complexes of edges with different orientations.

As we saw in Mishkin and Ungerleider's two visual systems hypothesis in Section 3.1, there are two distinct neural pathways for information leaving area V2. The ventral pathway and the dorsal pathway are each believed to be responsible for different types of visual processing.

- The *ventral pathway*, thought to be responsible for object identification and recognition, goes from V2 to V4 and then onward to areas in the inferior temporal cortex.
- The *dorsal pathway*, more focused on representations relevant to action, goes from V2 to areas V5 and V6 (also known as the middle temporal [MT] and dorsomedial [DM] areas).

We are currently interested primarily in the ventral pathway, because we are interested in how the visual system solves the representation learning problem of mapping unstructured stimuli in the visual field onto representations of objects in the distal environment – representations that can serve as inputs for more complicated forms of classification and information processing.

Area V4 takes the representations from V2, which contain primitive representations of shapes and some information about depth, and then uses those representations to



**Figure 12.6** An illustration of hierarchical visual processing. Stimuli are processed in a series of visual areas, which process features of increasing complexity. (Figure 1 from Manassi, Sayim, and Herzog 2013)

construct further representations that incorporate more information about figure/ground segmentation, as well as about colors. Area V4 is thought to be the first stage in early visual processing that is modulated by attention. Representations leaving V4 then go to the inferior temporal cortex (ITC), where there are various areas specialized for sensitivity to high-level patterns. Areas in the ITC include the fusiform face area (FFA), thought to be specialized for face recognition, and the fusiform body area (FBA), believed to be dedicated to identifying the human body and body parts.

The general pattern in the visual cortex is increasing determinacy of representation, with different areas in each pathway using the outputs of earlier areas to identify increasingly complex features and complexes of features. You can think about the different areas in the ventral pathway, for example, as forming a hierarchy. The information-processing hierarchy built into the ventral hierarchy is illustrated in Figure 12.6.

The reason for looking at the visual cortex in such detail is that this type of hierarchical information-processing is (by design) exactly what we find in many deep learning systems. Deep learning often involves multiple layers of information processing, each building on the outputs of earlier layers.

The general reasoning behind this approach is explained clearly in LeCun, Bengio, and Hinton 2015:

Deep neural networks exploit the property that many natural signals are compositional hierarchies, in which higher-level features are obtained by composing lower-level ones. In images, local combinations of edges form motifs, motifs assemble into parts, and parts form objects. Similar hierarchies exist in speech and text from sounds to phonemes, syllables, words and sentences. The pooling allows representations to vary very little when elements in the previous layer vary in position and appearance.

If, as the authors point out, many natural signals are hierarchically organized, then a deep learning system has to learn two things. First, it has to learn how to extract the bottom-level features. And then, second, it has to learn how to compound them to yield an accurate representation of the input in much the same way as the mammalian visual system has learned through evolution to parse the light energy hitting the retina into representations of objects.



Perhaps the principal challenge in doing this is to solve what LeCun, Bengio, and Hinton call the *selectivity/invariance problem*, which they illustrate with a canine example. Samoyeds are beautiful white dogs from the husky family. Imagine that you are trying to construct a system that will identify Samoyeds and distinguish them from other animals. Since Samoyeds have long fur, it is easy enough to distinguish them from short-haired dogs such as greyhounds. But what about distinguishing Samoyeds from their distant relative, the wolf?

The system needs to develop a highly *selective* representation of a Samoyed, one that will pick up on the relatively small differences between a Samoyed and a white wolf. But at the same time, the representation needs to be *invariant* in a way that ignores the conspicuous differences between what the animal looks like from different angles and in different postures. To illustrate the problem, think about how similar a Samoyed and a white wolf would look if they were both sitting down and viewed face on, and compare it to how different two Samoyeds would look if one was lying down and viewed from above, while the other was running and viewed from in front.

Deep neural networks are able to solve the selectivity/invariance problem because they are constructed as multilayered stacks of simple network. Each network performs a mapping from input to output, with the output from a given network serving as the input to the next network in the stack (just as the output of area V1 in the visual cortex is the input to area V2). Each of those mappings codes the input in terms of increasingly complex features (just as area V4 takes the primitive shapes in the output from area V2 and segments them into figure and ground to enhance the representation of depth).

The details of how this process works are extremely technical, but there are two fundamental types of neural network that feature in many deep learning systems and that can be explained in a relatively nontechnical way. The first are called *autoencoders*. The second are called *convolutional neural networks*. We turn to them in the next section.

## 12.3

### The Machinery of Deep Learning

A typical deep learning system is a multilayered stack of neural networks. In many respects, these neural networks are similar to many of the neural networks that we looked at in Chapter 5 and then again in Chapters 10 and 11. For example,

- they are often *feedforward*; that is, activation flows through the network in one direction – it does not loop back in feedback loops
- they contain *hidden layers*, in which the real information-processing work of the network gets done
- and they often learn through forms of *backpropagation*, with an error signal at the output from the network used to adjust the weights of units in the hidden layers

But the individual networks within deep neural networks also incorporate several features that we have not yet looked at.

In this section, we'll start by looking at a type of neural network called an *autoencoder*. Autoencoders are a very nice example of how a network can learn to extract features from unstructured raw data. Then we turn to what was probably the most important innovation in the deep learning revolution, namely, the development of convolutional neural networks (ConvNets). ConvNets are a special kind of feedforward neural network, particularly suited to processing data organized in a grid-like format (a photograph, for example, or the human visual field).

## Autoencoders

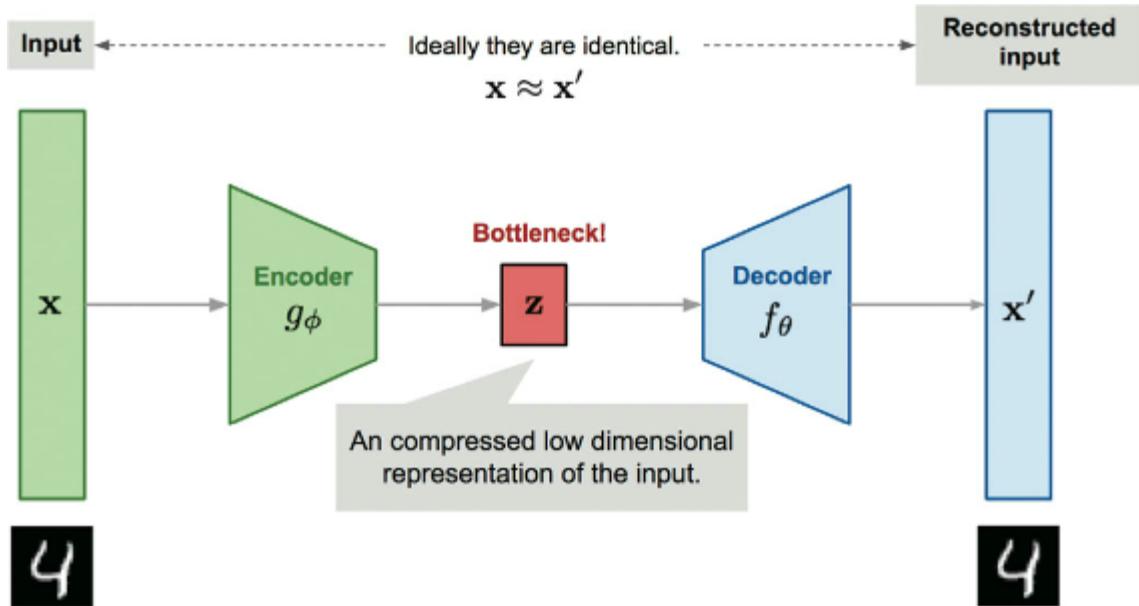
Think about the party game of charades. When you play charades, your job is to take a word or phrase and then mime it out so that the other players can guess what you started with. In other words, you take a signal as input and then represent it in such a way that the original signal can be recovered as output. Your mime performance is the representation and the other players' guess is the output. You have been successful if the output matches the input (i.e., if the other players come up with the word you started with).

This is pretty much what an autoencoder does. It is a single neural network that learns to match output to input. This might seem a very trivial thing to aim for – surely all the network has to do is repeat the input, which would guarantee success? But the point of an autoencoder is that it has constraints built into it that prevent it from solving the problem so simply – just as the rule in charades that you are only allowed to mime prevents you from simply telling the other players the word or phrase you started with. These constraints effectively force the network to find a simplified representation of the input. The representation needs to be simplified enough to satisfy the constraints built into the network, but not so simplified that the network cannot recover the original signal.

To put this all a little more technically, an autoencoder is a feedforward neural network whose job it is to approximate the identity function. It has two components. One component *encodes* the original signal, while the other *decodes* it. Both encoding and decoding are done by hidden layers. The aim of the network is to minimize the difference between input and output, and it is trained through some form of error minimization algorithm such as backpropagation. An autoencoder will typically have built into it some type of bottleneck that forces the network to find a more compressed way of representing the input.

So, for example, if the network has 100 input units and 100 output units, then the bottleneck might be a hidden layer with only fifty units. This constraint forces the network to compress an input signal that has 100 dimensions of variation into a representation that only allows variation in 50 dimensions. The network has to find key features of the input signal that will allow the encoded representation to be decoded back to something close to the original signal. Looking at it like this explains why feature learning is sometimes called *dimensionality reduction*.

There are two illustrations of autoencoders in Figure 12.7. The left-hand diagram is a schematic wiring diagram for an autoencoder neural network, while the second shows what such a network might actually do.



**Figure 12.7** Illustration of how an autoencoder compresses and then decompresses a signal. (Reproduced by permission of Lilian Wengweng)

Autoencoders can be used to *pretrain* deep neural networks, as an alternative both to starting with a random assignment of weights and to using some form of supervised learning process (where the network is given explicit feedback). Here is how it works. The autoencoder learns to compress its input signal in such a way that that signal can be decoded fairly accurately. This is a form of feature learning. What the network is doing is picking up on salient features in the input and using those features to encode the original signal. Once the network has learned to do this, it will encode any relevantly similar input signal in the same way. So, now the encoding component of the autoencoder can be built into a deep neural network as one of its layers.

To illustrate this, we can go back to our visual system example. Imagine that we are trying to build a deep neural network that approximates visual processing in the ventral stream (not as fanciful as it sounds – as we'll see shortly, this has actually been done). So, we need to build into it a network that will carry out edge detection. One possibility would be to train an autoencoder on images relevantly similar to those serving as inputs to area V1.

As the autoencoder learns to compress and decompress those images, it effectively learns to encode images as complexes of edges. We can then bypass the decoding component of the autoencoder and wire the bottleneck hidden layer up so that it becomes the input into the layer of the deep neural network corresponding to area V2. And then we might repeat the same process with an autoencoder that would learn to carry out V2-type feature detection, developing sensitivity to multiedge features.



### Exercise 12.3 Explain in your own words the basic idea behind an autoencoder.

In practice, there are better and more effective ways of constructing deep and multilayer, neural networks (e.g., by using the convolutional networks that we will look at next). But autoencoders are a good illustration of some basic deep learning principles.

## Convolutional Neural Networks

For a long time, computer vision was a final frontier for AI. It seemed almost impossible to design an artificial system to perform even a tiny fraction of the complicated visual information processing that a human child or adult performs almost every waking second of their lives. That has all changed in the last few years. The capabilities of the best modern machines are now not far short of their human models. This in turn is making possible a tsunami of technological advances, from facial recognition software in cellphones to self-driving cars. Convolutional neural networks (ConvNets) are playing a key role in this technological transformation.

The power of ConvNets was first revealed at a machine vision competition in 2012. The ImageNet Large-Scale Visual Recognition Challenge uses the massive ImageNet database created by Li Fei-Fei and Kai Li, on the faculty at Stanford and Princeton, respectively. ImageNet is actually derived from a lexical database called WordNet, which classifies English words and expressions into classes of synonyms called *synsets*. The goal of WordNet is really to provide an exhaustive catalog of all the concepts expressible in English. There are around 100,000 synsets in WordNet (around 80 percent associated with nouns and noun phrases) and ImageNet contains on average 1,000 images for each synset – and so contains around 10 million images in total.

Computer vision algorithms entered into the ImageNet competition had two different tasks to perform on a test database of images. The first task was to decide, for an evaluation set of image-category pairs, whether each image actually contained an example of an object from the relevant category (whether the image contained a washing machine, for example). The second task was to identify a specific object within an image and place a box around it (e.g., to identify a screwdriver centered on a specific pixel in the image).

The ImageNet competition started in 2010. Two years later, in 2012, a team from the University of Toronto entered a program called SuperVision. The team was led by Geoffrey Hinton who had been an early pioneer of the backpropagation learning for neural networks and then in the early 2000s was one of the researchers at CIFAR (the Canadian Institute for Advanced Research) responsible for the development of deep learning.

SuperVision was wildly successful. Its error rate on the test data set was only 16.4 percent. In comparison, the programs that won the competition in 2010 and 2011 had error rates of 28.2 percent and 25.8 percent, respectively, and the second-place program in 2012 had an error rate of 26.2 percent.

SuperVision's secret? As you've probably guessed, it was one of the first computer vision programs to use a convolutional neural network. The success of ConvNets in object



recognition and related tasks has been so compelling that almost all modern computer vision programs use them. And in fact, it did not take long to use ConvNets to develop programs much more successful even than SuperVision. The winner of the 2014 ImageNet competition, a program created by engineers from Google called GoogLeNet, had an error rate of only 6.7 percent.

ConvNets use a different type of mathematics from traditional neural networks. In traditional neural networks, activation spreads through the network via a form of matrix multiplication. This works because you can represent the activation level of neurons in each layer and the network weights as arrays of numbers (matrices). We can think of activation spreading through a neural network as a function of the activation levels of individual units and the weights that those individual units have. Matrix multiplication is the mathematical tool for combining arrays of numbers.

Convolution is a more complex mathematical operation. Basically, it creates a type of filter that functions as a localized feature detector. Imagine a grid composed of 10,000 pixels (organized in a square with dimensions 100 × 100). A convolutional layer will successively examine subgrids of this grid to filter out everything except the particular feature that it has been trained to detect. You can think of the process as a little square window (of size, say, 5 × 5) scanning the grid to detect edges, for example.

A ConvNet is a network that has at least one convolutional layer. According to *Deep Learning*, an influential textbook written by Ian Goodfellow, Yoshua Bengio, and Aaron Courville, ConvNets are designed with three characteristic features that together make them highly suited to processing data that has a grid-like organization (such as an image for example). Here they are (in slight different terminology):

- sparse connectivity
- shared weights
- invariance under translation

We will look at these in order.

## Sparse Connectivity

Like many traditional neural networks, ConvNets are feedforward. So, activation flows through the network, without any feedback loops. And they have multiple hidden layers. But the wiring is crucially different. In a traditional neural network, each unit in a given layer is connected to all the units in the previous layer and to all the units in the following layer. In other words, each unit is *fully connected*.

In contrast, the hidden layers in ConvNets are typically not fully connected. Each unit is wired to receive inputs only from a small number of units in the preceding layer, and it projects forward only to a small number of units in the next layer. This means that units are in effect specialized to respond to local regions of the input to any given layer. So, thinking about this general idea as it might apply to machine vision, you can think of individual units as specialized to respond to particular regions of the visual array.

The sparse connectivity feature of ConvNets is directly inspired by the mammalian visual system. As we saw earlier, information-processing in the early visual system often takes as input retinotopic maps that are more or less isomorphic to the patterns of light energy reaching the retina. Individual neurons in, say, area V1 have a receptive field that only covers a small part of that retinotopic map. By analogy, you can think of a given unit's receptive field as the units that project to it.

## *Shared Weights*

In ConvNets, as in traditional neural networks, every unit has a weight and learning takes place through changes in the weights. A unit's weight is really a measure of how much it contributes to the overall behavior of the network – a low weight indicates a small contribution and increasing the weight increases the significance of the unit to the network.

Characteristically, all the weights in a given layer of a traditional network are completely independent of each other. Each unit learns on its own, as it were. The backpropagation learning algorithm, for example, works by assigning to each individual unit a share in the responsibility for the overall degree of error. It then adjusts the weights of each unit to reduce the error attributable to that unit, which has the result of reducing the error in the network as a whole.

In ConvNets, on the other hand, there are dependencies between the weights of different units. A ConvNet might have multiple units with the same weight. Or it might have units whose weights are a function of the weights of other units. This feature adds further to the efficiency of ConvNets relative to traditional networks. Sparse connectivity means that ConvNets have fewer connections. Fewer connections means fewer computations as activation flows through the network. Reducing the number of different weights involved reduces the overall number of computations still further.

## *Invariance under Translation*

The last distinguishing feature of ConvNets arises from the first two, and in particular from how the shared weights are set up. To understand it, think again about hierarchical information processing in the visual system. We saw that the various areas in the ventral stream generate a succession of map-like representations of the distal environment, with each map building on earlier maps and adding more complex features. Simplifying somewhat, each area takes an input map, analyzes it, and then produces a richer output map. This process works because the structure of the original map is preserved through a series of successive transformations.

Now, the visual system often has to process retinal images derived from similar scenes – or from different perspectives on the same scene. It would be desirable for similarities in structure to be preserved across different episodes of vision. For an example, consider how



the object recognition systems in the inferior temporal cortex might recognize that the object you are now looking at is the same object you were looking at from a different angle a couple of minutes ago. The object recognition systems need to pick up on both similarities and differences. It needs to be sensitive to similarities to register that this is indeed the same object, but at the same time it needs to be sensitive to differences in the object's visual appearance due to the changed angle.

For the object recognition systems to be sensitive to similarities and differences in this way, it must be the case that similar sorts of changes in objects produce similar sorts of changes in the images that represent them. So, for example, whenever an object moves a certain distance to the left, the image of the object will move a certain distance to the left in the map that the visual system is building of the distal environment. This is the property of invariance under translation. ConvNets have the property of invariance under translation because changes in the input are systematically reflected in the output.



#### **Exercise 12.4** Write down brief definitions in your own words of the three characteristic features of ConvNets.

To conclude, then, ConvNets incorporate a range of features that make them particularly well suited to processing information that comes in a grid-like format. That is why they are so successful at object recognition and why they now dominate machine vision. But many types of information can be presented in a grid-like format. Linguistic information is an example. A sentence, after all, is really just a 1D grid of words. And so, it is not surprising that ConvNets have also been successfully applied in natural language processing. They have proved surprisingly good at speech recognition and semantic parsing, even though they have no linguistic information coded into them.

## 12.4

## Deep Reinforcement Learning

One of the main themes emerging from our discussion of deep learning is how biologically inspired it is. On a large-scale, deep neural networks are typically constructed in multiple layers to implement hierarchical information-processing explicitly modeled on information-processing in the primate visual system. On a smaller scale, the sparse connectivity and shared weights that we find in convolutional neural networks is much more neurally plausible than the full connectivity and atomistic learning that we find in traditional neural networks. We turn now to another example of biologically plausible deep learning.

As was extensively covered in the national and international press, the AlphaGo program created by Google's Deep Mind research group has been spectacularly successful at beating the world's leading experts at the game of Go. At the Future of Go summit in 2017, AlphaGo defeated Ke Jie, the number 1 player at the time, building on previous victories in 2016 (watched by tens of millions of people worldwide) over Lee Sedol, also one of a handful of top human Go players. These victories were widely recognized as historic achievements for



**Figure 12.8** A move in the Google DeepMind challenge between AlphaGo and Lee Sedol in 2016. (Handout/Getty Images)

AI. Go is even more complex than chess, since games are longer and there are many more legally permissible moves available (250 for Go, as opposed to 35 for chess).

AlphaGo was initially trained on a massive database of 30 million moves from an online server. The training was a form of supervised learning, with AlphaGo receiving explicit feedback on how successful it was. Once AlphaGo had achieved a relatively high level of playing strength, the training shifted to reinforcement learning. This is what we will focus on.

Reinforcement learning is not like any of the types of machine learning we have looked at up to now. Reinforcement learning differs from supervised learning. The network does not receive explicit feedback on how well it is doing, and it is not learning from a labeled training set. But nor is reinforcement learning a form of unsupervised learning, which is what we find in convolutional neural networks. In unsupervised learning networks learn to detect patterns in data. An unsupervised network might learn to become an edge detector, for example.

Reinforcement learning is distinctive because (unlike unsupervised learning) it *does* depend upon a feedback signal. But at the same time (unlike supervised learning) the feedback signal does not tell the network exactly what it has done wrong. Instead, reinforcement learning is driven by a reward signal. The job of a network incorporating a reinforcement learning algorithm is to maximize the reward. But it is not told how to do that. It has to work out for itself which outputs are most profitable, so that it can repeat and/or adapt them.

Here is a way of thinking about the feedback that drives reinforcement learning. The network is not receiving instructions. It is not being trained on typical exemplars of a classification task, for example. Nor is it being told how accurately it is performing the task. Instead its output is being evaluated through the reward signal. Once it has received the



evaluation/reward, it has to learn how to adapt the output to increase the reward. In a slogan, *supervised learning proceeds through instruction, while reinforcement learning proceeds through reward.*

Because the feedback in reinforcement learning is indirect, the network will have to engage in a degree of trial and error. The network will not increase its reward simply by repeating what has worked in the past. It needs to try strategies that it has not tried before in order to see whether they work better than current ones. So, successful reinforcement learning involves striking a balance between exploring new strategies and sticking to tried and tested strategies.

Reinforcement learning in AlphaGo was achieved by getting the network to play games of Go against former versions of itself (in fact, these former versions of itself were selected at random from earlier iterations of the supervised learning process). The reward signal was very simple: +1 for a win and -1 for a loss. AlphaGo then used a learning algorithm called *stochastic gradient descent*. Stochastic gradient descent works to reduce error (or, as it is often put, to minimize a cost function), where the error is the negative reward that comes from losing. The algorithm changes the weights in the network so that the negative reward eventually becomes a positive reward.

So, AlphaGo's training was a mixture of supervised learning and reinforcement learning. Human participation was crucial to the supervised learning phase, and then dropped out of the picture for the reinforcement learning phase. You might be wondering, then, whether it is possible to cut the human piece out of the equation completely. Could a network be trained from scratch to play Go, using only reinforcement learning and without any access to databases of position and games, or any kind of expert knowledge?

Well, a few months after officially retiring AlphaGo, the team at DeepMind announced a new version of the program, which they called AlphaGo Zero. As its name suggests, AlphaGo Zero incorporates zero supervised learning. The only thing it knows about Go when it starts is the rules, and so it starts off playing completely at random. Using an updated reinforcement learning algorithm, the network modifies its weights by playing against itself.

It turns out that AlphaGo Zero learned very quickly to beat previous iterations of AlphaGo. After three days (and 4.9 million games played against itself), it was able to beat the version of AlphaGo that had defeated Lee Sedol in 2016 – and it was a decisive victory, by 100 games to 0. Within forty days it was playing Go well enough to beat all existing versions of AlphaGo. Moreover, it was computationally much more efficient than AlphaGo, consuming only a fraction of the power and using only four of Google's proprietary TPU (tensor processing unit) chips.

Reinforcement learning is clearly a very powerful tool for machine learning. We can close this discussion by observing two very important respects in which reinforcement learning is significant from a broader AI perspective. AI is not just about building powerful programs and algorithms that will harness the power of supercomputers to solve difficult problems. AI is also about developing techniques that are recognizably intelligent. And an important index of intelligence in computer programs is that they solve problems in ways similar to how we humans solve them (as opposed to beating them into submission by sheer computational power).

So, with that in mind, we need to observe that reinforcement learning is the most widespread and powerful form of learning in the animal kingdom. Behaviorists may not be correct that all learning is ultimately reinforcement learning (or, in a different terminology, instrumental conditioning). But there is no doubt that many, if not most, intelligent behaviors in humans and animals are the result of an exquisite sensitivity to action-reward contingencies. By incorporating reinforcement learning, therefore, deep neural networks can make a strong claim to being intelligent problem-solvers.

And in fact, the biological plausibility of reinforcement learning algorithms in deep neural networks goes even deeper. There is intriguing evidence from neuroscience that important learning mechanisms in the brain incorporate elements that we also find in reinforcement learning algorithms. The key findings here come from studies of the neurotransmitter *dopamine* (a neurotransmitter is a chemical released by neurons to send signals to other neurons). Dopamine is known to be deeply implicated in reward processing (and for that reason has been much studied in the context of addiction).

Pioneering studies of dopaminergic neurons during the 1990s led to the *reward prediction error hypothesis* of dopamine activity. According to this hypothesis, first proposed by Read Montague, Peter Dayan, and Terrence Sejnowski in 1996, the job of dopaminergic neurons is to fire when the actual reward differs from the expected reward. In other words, to tie this back to machine learning, dopaminergic neurons provide a reward-based error signal, which is exactly what is needed for reinforcement learning. Moreover, although the details are too technical to go into here, it turns out that there are striking parallels between the firing behavior of dopaminergic neurons and how the reward signals are processed in widely used gradient descent learning algorithms.



## Summary

Section 12.1 showed how machine learning algorithms can construct expert systems that can in turn reproduce (and in some cases improve upon) the performance of human experts in tasks such as medical diagnosis and processing mortgage applications. We looked at a specific example in the ID3 algorithm. In Section 12.2 we looked at the limits of traditional machine learning, and in particular its dependence upon lengthy and complex processes of feature engineering to label and organize databases. Representation learning is the subfield of machine learning dedicated to constructing algorithms that will carry out feature engineering on raw data. The most successful examples of representation learning have come from deep learning neural networks, which we looked at in Sections 12.2, 12.3, and 12.4.

Deep learning networks are hierarchically organized, on the model of the mammalian visual cortex. A deep learning network is made up of layers of individual neural networks, each responsible for extracting more complex features from the original raw data. We looked at two examples of neural networks that can feature in layers of deep learning networks – autoencoders and convolutional neural networks. In Section 12.4 we looked at a further type of learning – reinforcement learning, which is neither supervised learning nor unsupervised. We saw how



reinforcement learning algorithms made possible two of the most spectacular examples of deep learning – the AlphaGo and AlphaGo Zero programs created by Google’s Deep Mind research team.

## Checklist

### Expert Systems and Machine Learning

- (1) Expert systems are designed to reproduce the performance of human experts in particular domains (e.g., medical diagnosis or financial services).
- (2) Expert systems typically employ decision rules that can be represented in the form of a decision tree.
- (3) One problem studied in the field of machine learning is developing an algorithm for generating a decision tree from a complex database.
- (4) Generating a decision tree in this way is an example of Newell and Simon’s heuristic search hypothesis.

### The ID3 Machine Learning Algorithm

- (1) ID3 looks for regularities in a database of information that allow it to construct a decision tree.
- (2) The basic objects in the database are called *examples*. These examples can be classified in terms of their *attributes*. Each feature divides the examples up into two or more classes.
- (3) ID3 constructs a decision tree by assigning attributes to nodes. It assigns to each node the attribute that is most informative at that point.
- (4) Informativeness is calculated in terms of *information gain*, which is itself calculated in terms of *entropy*.

### Representation Learning and Deep Learning

- (1) Problem solving is often highly dependent upon how data sets are labeled through *feature engineering*, typically carried out by human experts.
- (2) Traditional machine learning algorithms such as ID3 work on databases that are already highly organized and labeled.
- (3) Representation learning is the subfield of machine learning dedicated to designing algorithms that will do their own feature engineering on raw data.

**Deep learning is a special kind of representation learning, carried out by multilayered artificial neural networks, hierarchically organized to extract increasingly complex features from input raw data.**

- (1) The hierarchical organization of deep learning networks is inspired by the hierarchical organization of the mammalian visual system.
- (2) An autoencoder is an example of a deep learning network that can learn to extract features from raw data.
- (3) Autoencoders compress raw data through a bottleneck so that it can subsequently be decoded.

- (4) Convolutional neural networks (ConvNets) are particularly well suited to feature engineering in data presented in a grid-like format.
- (5) ConvNets have been spectacularly successful in machine vision.
- (6) ConvNets feature (a) sparse connectivity, (b) shared weights, (c) invariance under translation.

**Reinforcement learning is a type of machine learning distinct both from supervised learning (in which networks are given explicit feedback on what they have done wrong) and unsupervised learning (such as representation learning).**

- (1) In reinforcement learning networks receive a reward signal, rather than an explicit error signal.
- (2) Reinforcement learning algorithms aim to maximize expected reward, and so they need to engage in trial and error to discover new reward-generating strategies.
- (3) Reinforcement learning is key to the two Go-playing networks, AlphaGo and AlphaGo Zero, developed by Google's Deep Mind research team.
- (4) Reinforcement learning is biologically plausible, both because reinforcement learning is the most common type of learning in the animal kingdom, and because of intriguing parallels between the firing activity of dopaminergic neurons and how reward signals are processed in widely used learning algorithms.

## Further Reading

Much of the literature in traditional machine learning is very technical, but there are some accessible introductions. Haugeland 1985 and Franklin 1995 remain excellent introductions to the early years of AI research. Russell and Norvig 2009 is much more up to date. Also see Poole and Mackworth 2010, Warwick 2012, and Proudfoot and Copeland's chapter on artificial intelligence in *The Oxford Handbook of Philosophy of Cognitive Science* (Margolis, Samuels, and Stich 2012). Medsker and Schulte 2003 is a brief introduction to expert systems, while Jackson 1998 is one of the standard textbooks. The *Encyclopedia of Cognitive Science* also has an entry on expert systems (Nadel 2005). See the online resources for a very useful collection of machine learning resources.

The application of ID3 to soybean diseases described in Section 12.2 was originally reported in Michalski and Chilausky 1980. The database for the tennis example explored in Section 12.1 comes from chapter 3 of Mitchell 1997. Wu et al. 2008 describes more recent extensions of ID3, including C4.5 and C5.0, as well as other data mining methods.

LeCun, Bengio, and Hinton 2015 is an article-length introduction to deep learning and the source of the definition of deep learning in Section 12.2. For more detailed discussion of the theoretical background, see the textbook *Deep Learning* (Goodfellow, Bengio, and Courville 2016), which is also available online for free by agreement with the publisher at [www.deeplearningbook.org](http://www.deeplearningbook.org). For a direct deep learning model of the visual cortex, see Cadieu et al. 2014. For the SuperVision program entered in the ImageNet competition in 2012, see Krizhevsky, Sutskever, and Hinton 2012.



The Deep Mind research group at Google has a website at [www.deepmind.com](http://www.deepmind.com) and has published articles in *Nature* on AlphaGo (Silver et al. 2016) and AlphaGo Zero (Silver et al. 2017). For a more general discussion of reinforcement learning, the classic textbook is Sutton and Barto 1998, with a second edition currently in preparation and available online in draft form from Richard Sutton's website at [www.incompleteideas.net](http://www.incompleteideas.net). Chapter 15 of the second edition covers the neuroscience dimension of reinforcement learning. The pioneering paper on the reward prediction error hypothesis was Montague, Dayan, and Sejnowski 1996. For a survey of the theory and experimental support, see Glimcher 2011.





## CHAPTER THIRTEEN

# Exploring Mindreading

### OVERVIEW 335

#### 13.1 Pretend Play and

#### Metarepresentation 336

The Significance of Pretend Play 336

Leslie on Pretend Play and

Metarepresentation 337

The Link to Mindreading 341

#### 13.2 Metarepresentation, Autism, and

#### Theory of Mind 341

Using the False Belief Task to Study

Mindreading 342

Interpreting the Results 344

Implicit and Explicit Understanding of False

Belief 347

#### 13.3 The Mindreading System 348

First Steps in Mindreading 349

From Dyadic to Triadic Interactions: Joint Visual

Attention 351

TESS and TOMM 352



## Overview

This chapter introduces what is often called *mindreading*. This is a very general label for the skills and abilities that allow us to make sense of other people and to coordinate our behavior with theirs. Our mindreading skills are fundamental to social understanding and social coordination.

The dominant model of mindreading in cognitive science emerged from studies of pretending in young children. Section 13.1 presents the information-processing model of pretense proposed by the developmental psychologist Alan Leslie. In Section 13.2 we build up from pretending to mindread, looking in particular at the *false belief task*, which tests young children's understanding that other people can have mistaken beliefs about the world.

The central feature of Leslie's model is what he calls the *theory of mind mechanism* (TOMM). The TOMM's job is to identify and reason about other people's *propositional attitudes* (complex mental states, such as beliefs, desires, hopes, and fears). Section 13.3 introduces a model of the entire mindreading system developed by the developmental psychologist and autism specialist Simon Baron-Cohen in response to a wide range of experimental data both from normal development and from autism and other pathologies.


**13.1**

## Pretend Play and Metarepresentation

Developmental psychologists think that the emergence of pretend play is a major milestone in cognitive and social development. Children start to engage in pretend play at a very young age, some as early as 13 months. Normal infants are capable of engaging in fairly sophisticated types of pretend play by the end of their second year. The evidence here is both anecdotal and experimental. Developmental psychologists such as Jean Piaget have carried out very detailed longitudinal studies of individual children over long periods of time. There have also been many experiments exploring infants' emerging capacities for pretend play.

The development of pretend play in infancy appears to follow a fairly standard trajectory. The most basic type is essentially *self-directed* – with the infant pretending to carry out some familiar activity. The infant might, for example, pretend to drink from an empty cup, or to eat from a spoon with nothing on it. The next stage is *other-directed*, with the infant pretending that some object has properties it doesn't have. An example of this might be the infant's pretending that a toy vehicle makes a sound, or that a doll is saying something. A more sophisticated form of pretense comes with what is sometimes called *object substitution*. This is when the infant pretends that some object is a different object and acts accordingly – pretends that a banana is a telephone, for example, and talks into it. Infants are also capable of pretense that involves imaginary objects. Imaginary friends are a well-known phenomenon.

Some forms of pretend play exploit the young infant's emerging linguistic abilities. Others exploit the infant's understanding of the different functions that objects can play. A common denominator in all instances of pretend play is that in some sense the infant is able to represent objects and properties not perceptible in the immediate environment – or at least, not perceptible in the object that is the focus of the pretense (since there may be a telephone elsewhere in the room, for example).

### The Significance of Pretend Play

Alan Leslie calls the infant's basic representations of the environment its *primary representations*. Primary representations include both what the infant perceives, and its stored knowledge of the world.

Leslie's model of infant pretense starts off from three basic observations:

- 1 *Pretend play in the infant depends crucially on how the infant represents the world (and hence on her primary representations).* If an infant pretends that a banana is a telephone, then she must be representing the banana to start with. The infant is in some sense taking her representation of a banana and making it do the job of a representation of a telephone. Similarly, the infant cannot represent a toy car as making a noise unless she is representing the car.
- 2 *We cannot explain what is going on in pretend play simply with reference to the infant's primary representations.* We cannot assume that the infant is somehow coordinating her banana



representation and her telephone representation. The problem is that the primary representation and the pretend representation typically contradict each other. After all, the banana is a banana, not a telephone.

- 3 *The pretend representations must preserve their ordinary meanings in pretend play.* During pretend play the infant cannot lose touch of the fact that, although she is pretending that it is a telephone, what she has in front of her is really a banana. Likewise, representing the banana as a telephone requires representing it as having the properties that telephones standardly have.

Putting these three ideas together, Leslie arrived at the idea that, although representations in pretend play have their usual meaning, they cannot actually be functioning as primary representations. They are somehow “quarantined” from ordinary primary representations. Without this sort of quarantining, the infant’s representations of the world would be completely chaotic. One and the same cup would be both empty and contain water, for example. The key problem is to explain how the quarantining takes place.

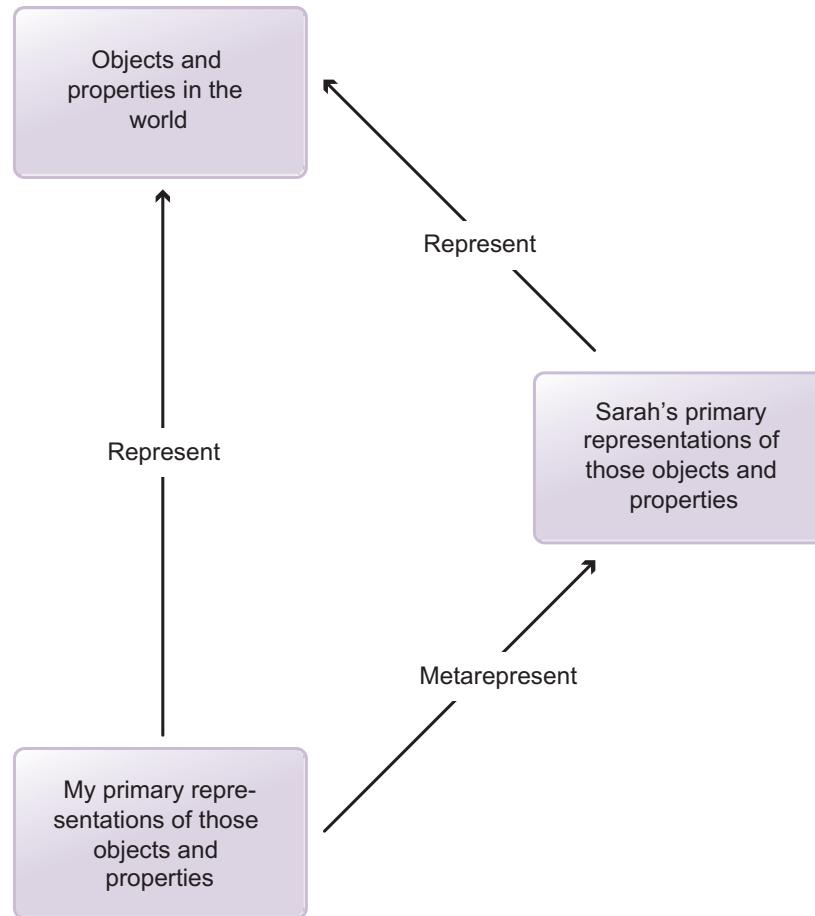
Leslie’s explanation rests upon a very basic parallel between how representations function in pretend play and how they function when we are representing other people’s mental states in mindreading. When we represent what other people believe or desire, we do so with representations that are also quarantined from the rest of our thinking about the world.

Suppose, for example, that I utter the sentence “Sarah believes that the world is flat.” I am asserting something about Sarah, namely, that she believes that the world is flat. But I am certainly not saying that the world is flat. If I were to utter the words “the world is flat” on their own, then I would standardly be making an assertion about the world. But when those very same words come prefixed by the phrase “Sarah believes that . . .” they function very differently. They are no longer being used to talk about the world. I am using them to talk about Sarah’s state of mind. They have become *decoupled* from their usual function.

## Leslie on Pretend Play and Metarepresentation

Let us look at this in more detail. When I describe Sarah as believing that the world is flat the phrase “the world is flat” is being used to describe how Sarah herself represents the world. Philosophers and psychologists typically describe this as a case of *metarepresentation*. Metarepresentation occurs when a representation is used to represent another representation, rather than to represent the world. The fact that there is metarepresentation going on changes how words and mental representations behave. They no longer refer directly to the world. But they still have their basic meaning – if they lost their basic meaning then they couldn’t do the job of capturing how someone else represents the world.

The basic picture is summarized in Figure 13.1. As the figure shows, my primary representations can serve two functions. They can represent the world *directly*. This is the standard, or default use. But they can also be used to metarepresent someone else’s primary representations. This is what goes on when we engage in mindreading.

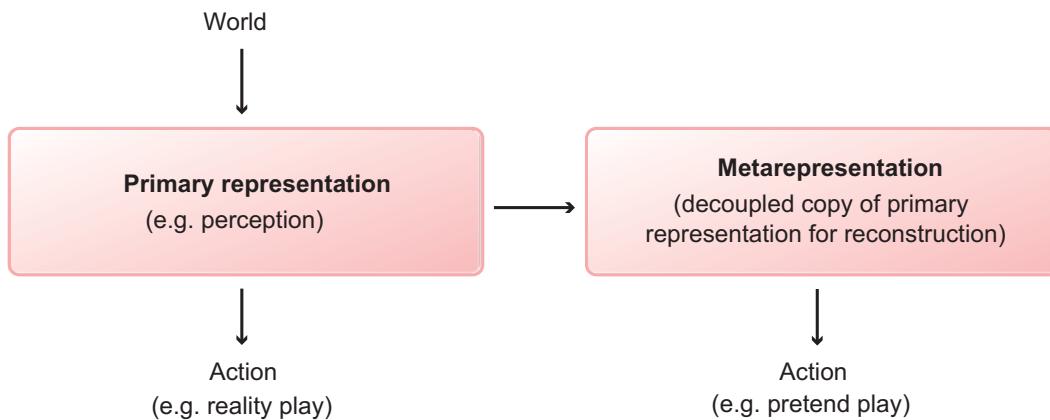


**Figure 13.1** An example of metarepresentation. Metarepresentation involves second-order representations of representations. In this example, I am representing Sarah's representations of certain objects and properties in the world.

Leslie's basic idea is that primary representations function exactly the same way in pretend play as when used to metarepresent someone else's state of mind. In both cases, primary representations are decoupled from their usual functions. In fact, Leslie argues, the mechanism that decouples primary representations from their usual functions is exactly the same in pretend play and in mindreading. The structure of Leslie's model is outlined in Figure 13.2.

Two features of the model deserve more discussion. First, the model incorporates a way of marking the fact that a primary representation has been decoupled and is now being used for pretend play. Second, it includes a way of representing the relation between agents and decoupled representations.

Leslie proposes that decoupling is achieved by a form of quotation device. In ordinary language we use quotation marks to indicate that words are being decoupled from their normal function. In fact, we often do this when we are reporting what other people have said. So, for example, the following two ways of reporting what Sarah said when she expressed her belief that the world is flat are more or less equivalent:



**Figure 13.2** The general outlines of Leslie's model of pretend play. (Adapted from Leslie 1987)

- (1) Sarah said that the world is flat.
- (2) Sarah said: "The world is flat."

The second report contains a device that makes explicit the decoupling that is achieved implicitly in the first report. His suggestion, then, is that the physical symbol system responsible for pretend play contains some sort of quotation device that can be attached to primary representations to mark that they are available for pretend play.



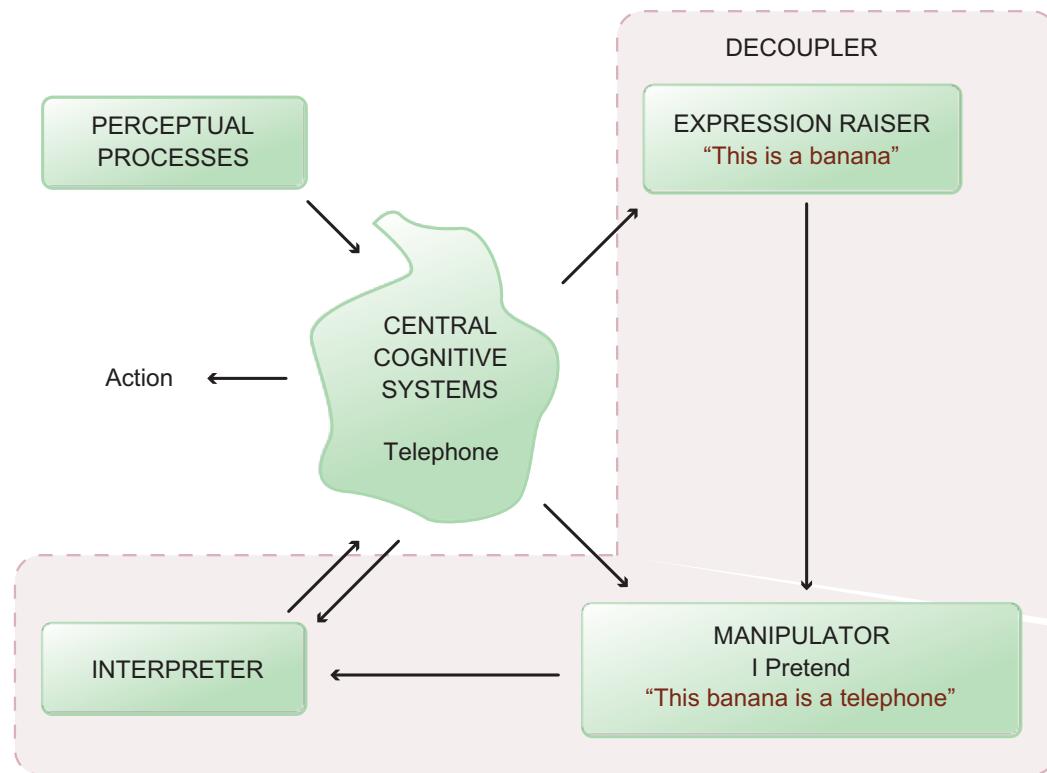
### Exercise 13.1 (1) and (2) are not completely equivalent. Explain why not.

How are decoupled primary representations processed in pretend play? Leslie's metarepresentation system contains a special operation, which he calls the PRETEND operation. The subject of the PRETEND operation is an agent (which may be the pretending infant himself). The PRETEND operation is applied to decoupled primary representations. But these are not pure decoupled representations. The essence of pretend play is the complex interplay between ordinary primary representations and decoupled primary representations. Leslie's model aims to capture this with the idea that decoupled representations are, as he puts it, *anchored* to parts of primary representations.

Let's go back to our example of the infant pretending that the banana is a telephone. The infant's representation of the banana is decoupled and then anchored to her primary representation of a telephone. Leslie would represent what is going on here in the following way:

I PRETEND "*This banana: it is a telephone.*"

The object of the PRETEND operation is the complex representation: "*This banana: it is a telephone.*" As the quotation marks indicate, the complex representation as a whole is decoupled. But it is made up of two representations – a (decoupled) representation of a banana and an ordinary representation of a telephone. The ordinary representation of the telephone is the anchor for the decoupled representation of the banana.



**Figure 13.3** Leslie's Decoupler model of pretense. This model makes explicit how the right-hand side of Figure 13.2 is supposed to work. (Adapted from Leslie 1987)

The details of Leslie's model of pretense can be seen in Figure 13.3. Here is how it works.

- Information goes from central cognitive systems into what Leslie calls the Expression Raiser. This is the system that decouples primary representations – by placing them within some analog of quotation marks.
- Decoupled primary representations can then be fed into the Manipulator, which applies the PRETEND operation as described earlier.
- The job of the Interpreter is to relate the output of the Manipulator to what the infant is currently perceiving. The Interpreter essentially controls the episode of pretend play. Pretend play requires certain inferences (for example – the inference that, since the telephone is ringing, I must answer it). These are implemented in the Interpreter, using general information about telephones stored in central systems.

Leslie's model explains both how infants can engage in pretend play, and how they can understand pretense in other people. In Figure 13.3 the infant herself is the agent of the PRETEND operation, but the agent could equally be someone else. This allows the infant to engage in collaborative pretend play – and, moreover, gives her an important tool for making sense of the people she is interacting with.



## The Link to Mindreading

Understanding that other people are pretending is itself a form of mindreading. In this sense, therefore, Leslie's model of pretense is already a model of mindreading. But, as Leslie himself points out, the basic mechanism of metarepresentation at the heart of the model can be applied much more widely to explain other forms of mindreading. This is because many forms of mindreading exploit decoupled representations, as we saw earlier. And so, once the basic mechanism of decoupling is in place, being able to perform other types of mindreading depends upon understanding the corresponding operations.

So, we might expect there to be operations BELIEVE, DESIRE, HOPE, FEAR, and so on, corresponding to the different types of mental state that a mindreader can identify in other people. These operations will all function in the same way as the PRETEND operation. At an abstract level these operations are all applied to decoupled representations. In order to represent an agent as believing a particular proposition (say, the proposition that it is raining), the mindreader needs to represent something of the following form:

Agent BELIEVES "It is raining."

where "it is raining" signifies a decoupled primary representation. This is exactly the same decoupled representation that would be exploited when the infant pretends that it is raining.

If this is right, then the foundations for the mindreading system are laid during the second year of infancy, when infants acquire the basic machinery of decoupling and metarepresentation. It is a long journey from acquiring this basic machinery to being able to mindread in the full sense. Mindreading is a very sophisticated ability that continues to develop throughout the first years of life. Many of the operations that are exploited in older children's and adults' mindreading systems are much harder to acquire than the PRETEND operation. This is what we'll look at now.

## 13.2

## Metarepresentation, Autism, and Theory of Mind

In developing his model Leslie placed considerable weight on studies of children with autism. Autism is a developmental disorder that has been increasingly discussed and studied in recent years. Autism typically emerges in toddlers and the symptoms are often detectable before the age of 2. The disorder is strongly suspected to be genetic in origin, although its genetic basis is not very well understood.

For psychologists and cognitive scientists, autism is a very interesting disorder because it typically involves deficits in social understanding, social coordination, and communication. But these social and communicative problems are not typically accompanied by general cognitive impairments. Autistic subjects can have very high IQs, for example. Their problems seem to be relatively circumscribed, although autistics often have sensory and motor problems, in addition to difficulties with language.

One feature of autism that particularly sparked Leslie's attention is that autistic children have well-documented problems with pretend play. This has been revealed by many studies showing that pretend play in autistic children is very impoverished, in comparison both with much younger normal children and with mentally retarded children of the same age. In fact, the phenomenon is so widespread in autism that it has become a standard diagnostic tool. Parents are often first alerted to autism in their children by their apparent inability to engage in pretense and make-believe – and by the child's inability to understand what other people are up to when they try to incorporate the child into pretend play. And one of the first questions that clinicians ask when parents suspect that their child has autism is whether the child engages in pretend play.

This well-documented fact about autistic children is particularly interesting in the context of the other problems that autistic children have. These problems cluster around the very set of abilities in social understanding and social coordination that we are collectively terming mindreading. In 1985 Leslie was one of the authors of a very influential paper arguing that autistic children had a very specific mindreading deficit – the other two authors were Simon Baron-Cohen and Uta Frith.

## Using the False Belief Task to Study Mindreading

Baron-Cohen, Leslie, and Frith studied three populations of children. The first group were autistic, aged between 6 and 16 (with a mean of 11;11 – i.e., 11 years and 11 months). The second group of children suffered from Down syndrome, which is a chromosomal disorder usually accompanied by mental disability, often severe. The Down syndrome children varied from 7 to 17 years old (with a mean of 10). The third group (the control group) were children with no cognitive or social disorders, aged from 3;5 to 6, with a mean of 4;5.

It is very interesting to look at the overall cognitive ability of the three different populations, as measured on standard tests of verbal and nonverbal mental age, such as the British Picture Vocabulary test (which measures the ability to match words to line drawings) and the Leiter International Performance Scale (which measures nonverbal abilities such as memory and visualization).

The normal children scored lowest on the nonverbal measures. The normal children's mean nonverbal mental age of 4;5 compared to a mean nonverbal mental age of 5;1 for the Down syndrome group and 9;3 for the autistic group. The Down syndrome group had the lowest verbal mental age (with a mean of 2;11). The verbal skills of the autistic group were significantly ahead of the normal children (with a mean verbal mental age of 5;5). These numbers are all depicted in Table 13.1.

Baron-Cohen, Leslie, and Frith tested the mindreading abilities of the three groups by using a very famous experimental paradigm known as the *false belief test*. The false belief test was first developed by the developmental psychologists Heinz Wimmer and Joseph Perner in an article published in 1983.

There are many different versions of the false belief test, but they all explore whether young children understand that someone might have mistaken beliefs about the world.

**TABLE 13.1** The three groups studied in Baron-Cohen, Leslie, and Frith (1985)

<b>POPULATION</b>	<b>MEAN VERBAL MENTAL AGE</b>	<b>MEAN NONVERBAL MENTAL AGE</b>
Normal group	4;5	4;5
Down syndrome group	2;11	5;1
Autistic group	5;5	9;3

There is a very basic contrast between belief, on the one hand, and knowledge, say, on the other. Consider knowledge. There is no way in which I can know that some state of affairs holds without that state of affairs actually holding. Knowledge is an example of what philosophers sometimes call *factive* states.



### Exercise 13.2 Can you give examples of other mental states that are factive in this sense?

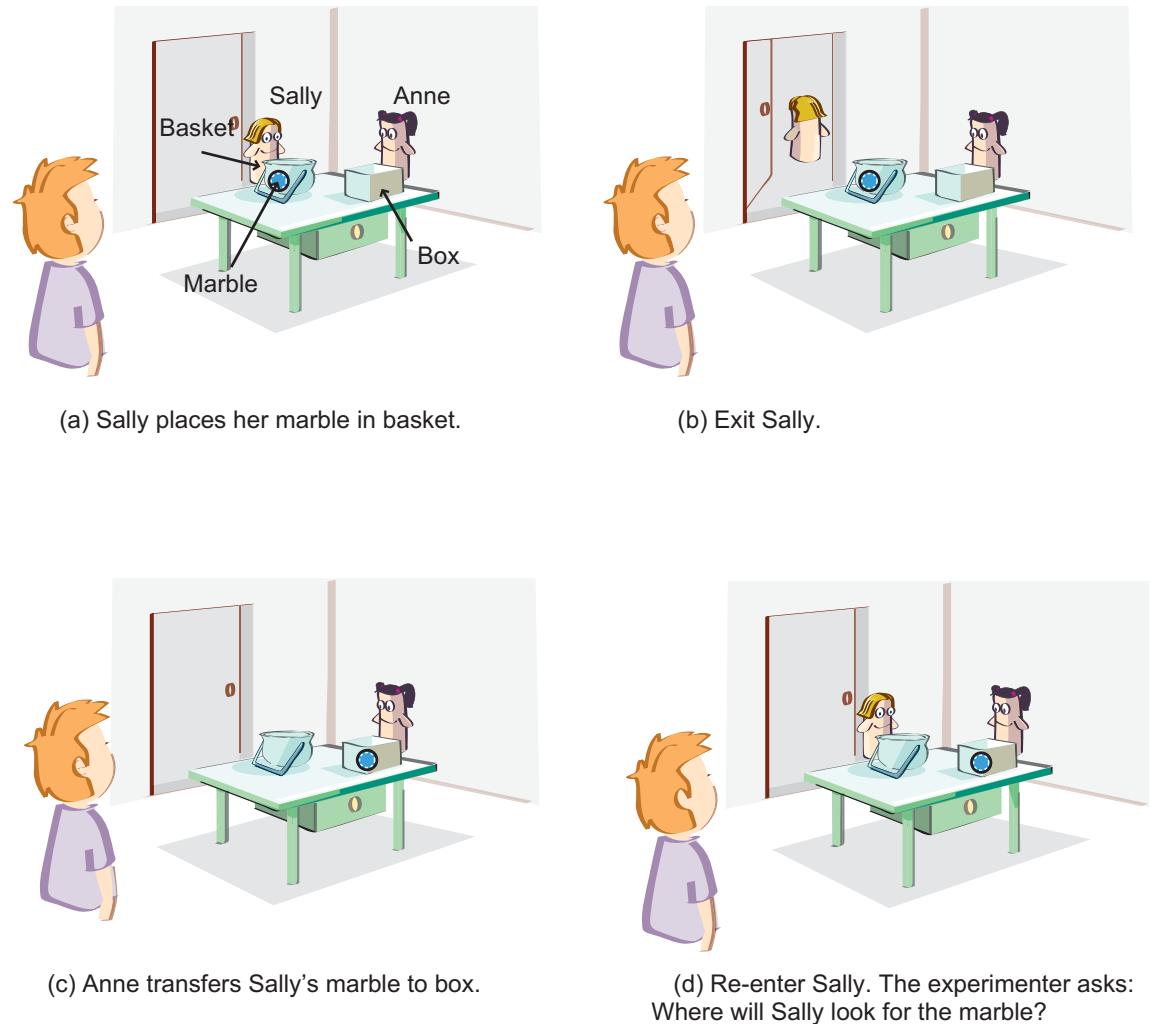
In contrast, beliefs are not factive. I cannot have false knowledge, but I can (all too easily) have false beliefs. And anyone who understands what belief is needs to understand that it is possible to have false beliefs.

So, if a young child does not understand the possibility that someone might have false beliefs about the world, then there seems to be no sense in which they understand what is involved in believing something. They cannot possess the concept of belief. And this, in turn, tells us something about their mindreading skills. Children who do not understand the concept of belief are lacking a fundamental component of the mindreading tool kit.

But how do we test whether children understand the possibility of false belief? This is where the false belief test comes into the picture.

The experimental setup used by Baron-Cohen, Leslie, and Frith is a variant of Wimmer and Perner's original false belief test. It is depicted in Figure 13.4. The child being tested is seated in front of an experimenter, who has two puppets, Sally and Anne. Between the child and the experimenter is a table with a basket and box. In front of the child, Sally places a marble in the basket and then leaves the room. While she is away Anne transfers the marble from the basket to the box. Sally then returns. The experimenter asks the child: "Where will Sally look for her marble?" (or, in some versions of the test, "Where does Sally think the marble is?").

The point of the experiment is that, although the child saw the marble being moved, Sally did not. So, if the child has a clear grip on the concept of belief and understands that it is possible to have false beliefs, then she will answer that Sally will look in the basket, since nothing has happened that will change Sally's belief that the marble is in the basket. If, on the other hand, the child fails to understand the possibility of false belief, then she will answer that Sally will look for the marble where it in fact is, namely, in the box.



**Figure 13.4** The task used by Baron-Cohen, Leslie, and Frith to test for children's understanding of false belief. (Adapted from Baron-Cohen, Leslie, and Frith 1985)



**Exercise 13.3** Explain in your own words the logic behind the false belief task. Do you think it succeeds in testing a young child's understanding of false belief?

## Interpreting the Results

The results of the experiment were very striking. The main question that the experimenters asked was the obvious one, which they called the Belief Question: "Where will Sally look for her marble?" But they also wanted to make sure that all the children understood what was going on. So they checked that each child knew which doll was which and asked two further questions:

- |  |                        |
|--|------------------------|
| "Where was the marble in the beginning?" | (the Memory Question)  |
| "Where is the marble really?"            | (the Reality Question) |



### Exercise 13.4 Explain in your own words the purpose of asking these two extra questions.

Baron-Cohen, Leslie, and Frith found that all the children understood the experimental scenario. None of them failed either the Memory Question or the Reality Question. But there was a very significant difference in how the three groups fared with the Belief Question. Both the Down syndrome group and the normal group were overwhelmingly successful – with correct answers from 86 percent and 85 percent, respectively. This is despite the fact that the Down syndrome group had a mean verbal mental age of less than 3. In very striking contrast, the autistic group (with a mean verbal mental age of 5;5) performed extremely poorly. In fact, 80 percent of the autistic children failed the Belief Question, despite a relatively high level of general intelligence.

The experimenters concluded that autistic children have a highly specific mindreading deficit. As they put it in the original paper, “Our results strongly support the hypothesis that autistic children as a group fail to employ a theory of mind. We wish to explain this failure as an inability to represent mental states. As a result of this the autistic subjects are unable to impute beliefs to others and are thus at a grave disadvantage when having to predict the behavior of other people” (Baron-Cohen et al. 1985: 43).

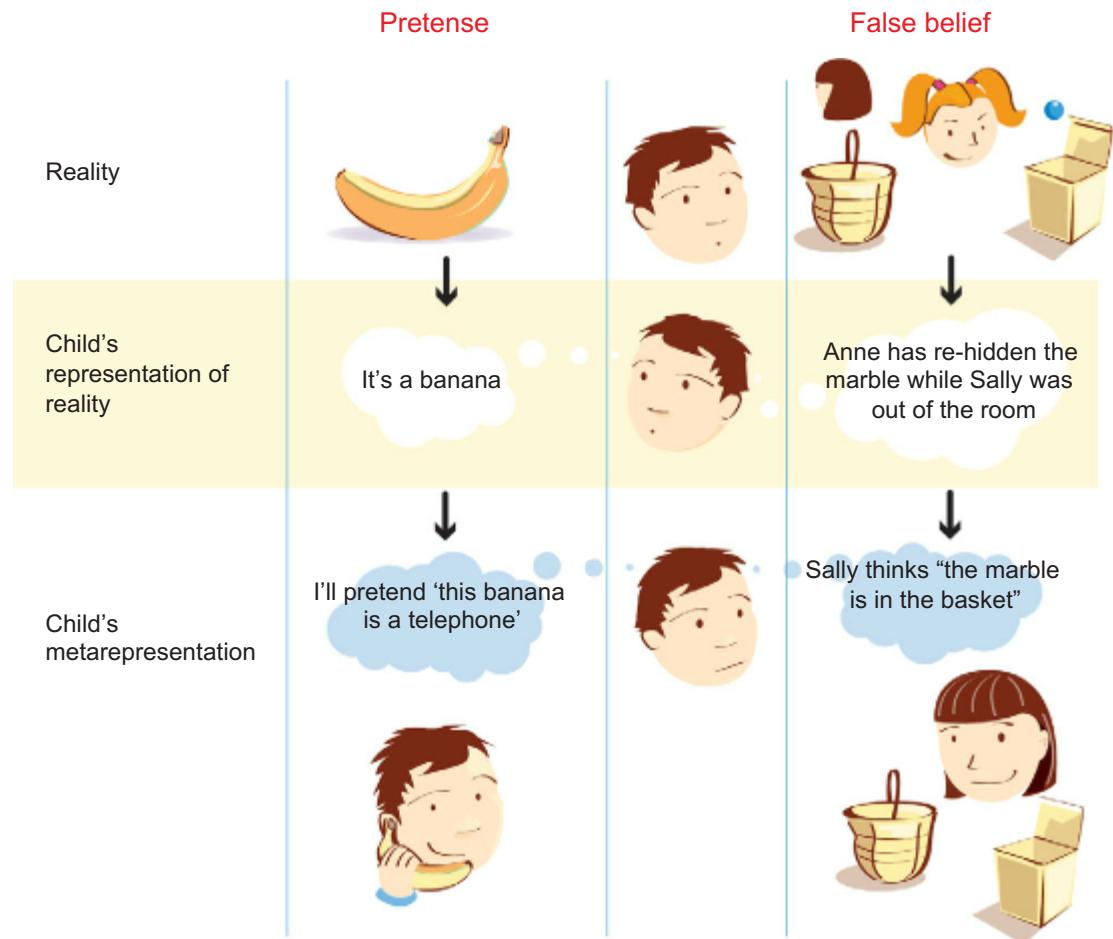
Notice the specific diagnosis of why the autistic children fail the false belief task. It is described as a failure in the ability to represent mental states – in metarepresentation. This connection with Leslie’s theory of pretend play is illustrated in Figure 13.5.

Leslie’s theory allows us to connect two things that seem on the face of it to be completely unconnected. The first is the fact that autistic children have severe problems with pretend play. The second is that autistic children have serious difficulties with the false belief task – and so, many have concluded, with mindreading more generally. These two things turn out to be very closely connected if we think that both pretend play and mindreading critically depend upon metarepresentation. Autistic children’s difficulties with pretend play and with mindreading turn out to have a common cause and a common explanation, namely, a deficit in metarepresentation.

This way of thinking about what is going wrong in the social development of the autistic child goes hand in hand with a model of how social development progresses for the normal child. On Leslie’s model, as reinforced by the experimental studies we have been examining, pretend play has a crucial role to play in the emergence of metarepresentation. In autistic children, for reasons that are not yet understood, the process of developing metarepresentational abilities never really gets going.

Here is the basic picture. Pretend play rests upon a basic portfolio of metarepresentational abilities. These metarepresentational abilities permit primary representations to be decoupled from their usual functions. Once decoupled they can serve as inputs to the PRETEND operation. The same basic machinery is supposed to be exploited in mindreading more generally.

When young children (or adults, for that matter) successfully pass the false belief task, they are (according to the model) starting with their memory of the ball being placed in the basket. The metarepresentational mechanisms allow this primary representation to be



**Figure 13.5** Illustration of the connection between pretend play and success on the false belief task.

decoupled from its usual role (so that, for example, it is not invalidated by watching Anne transfer the marble from the basket to the box). This allows the child to form a representation along these lines:

Sally BELIEVES “The marble is in the basket.”

There is still a very important gap in the account, however. The problem is chronological. Pretend play emerges during the second year of life. But children do not typically pass the false belief test until they are nearly 4. There is a very clear sense, therefore, in which the BELIEVES operation must be much harder to acquire than the PRETENDS operation. But why is this? And what is the developmental progression that takes the normal child from pretend play to successful mindreading, as evidenced by success on the false belief task? We will come back to these questions later. First, though, we need to consider some important experiments suggesting that children may be able to understand false beliefs significantly earlier than suggested by the standard false belief task.



## Implicit and Explicit Understanding of False Belief

The false belief task originally proposed by Baron-Cohen, Leslie, and Frith is a verbal task. Children are explicitly asked about where they think Sally will look, or where they think the marble is. But it may be that these explicit questions introduce additional computational demands that muddy the waters. Perhaps young children fail the false belief task because they cannot cope with these extra computational demands, rather than because they do not understand false belief?

One way of exploring this possibility would be to develop a less demanding false belief test. This was done by Kristine Onishi and Renée Baillargeon in a famous set of experiments first published in 2005. Instead of explicitly asking children about how the characters they were observing would behave, or what they believed, Onishi and Baillargeon used a violation of expectations paradigm that measured looking times. (Rather like the experiments on infant folk physics discussed in Chapter 11.)

In their experimental setup, 15-month-old infants were familiarized to an actor searching for a toy in one of two boxes (yellow and green, respectively). They were then presented with different conditions. In one condition the toy was moved from one box to the other with the actor clearly watching. In a second condition the toy was moved in the absence of the actor. After the toy was moved the actor then looked for the toy in one of the two baskets.

Onishi and Baillargeon hypothesized that the length of time that the infants looked at each of the scenarios would be a guide to their implicit understanding of false belief. Consider the second scenario, where the toy is moved without the actor seeing. Suppose that the toy was moved from the green box to the yellow box without the actor observing. Then the actor would presumably have a false belief about the toy's location, thinking it to still be in the green box when it is really in the yellow box. If infants understand this then they will presumably expect the actor to search in the green box. This expectation will be violated if the actor searches in the yellow box.

The experiments revealed a robust effect. Infants looked significantly longer when the actor searched in the yellow box than when the actor searched in the green box, even though the toy was really in the green box. Onishi and Baillargeon claim that the infants were surprised that the actor did not act on the basis of his (false) belief that the toy was still in the green box. So, they conclude, infants have an understanding of false belief much earlier than suggested by the traditional false belief task.

The Onishi and Baillargeon results have been replicated and expanded by other researchers. At the same time, however, there has been considerable debate about how to interpret them. Some cognitive scientists, including Onishi and Baillargeon themselves, think that the results show that young infants have a full understanding of false belief, directly refuting the standard claim that children do not arrive at a full understanding of false belief until around 4 years of age. Others take a more measured approach. This is what we shall do here.

The original, verbal false belief experiments seem to be testing for a cognitive ability considerably more sophisticated than could be revealed by the Onishi and Baillargeon experiments. The earlier experiments are directly targeting explicit conceptual abilities manifested in verbal responses and explicit reflection. Children are asked about what agents will do and what they believe. What the experiments are getting at is mastery of the concept of belief, together with the complicated vocabulary and other baggage that goes with it.

In contrast, the Onishi and Baillargeon experiments are probing the nonverbal expectations that young children have about behavior and how behavior is affected by what an agent has and has not observed. It is clear that these are related in at least one important sense. Nobody who lacked the nonverbal expectations identified in the Onishi and Baillargeon experiments could possibly pass the much more sophisticated false belief test. At the same time, though, the dependence doesn't seem to hold in the opposite direction. It seems perfectly possible to have the right nonverbal expectations without being able to articulate them in the right sort of way to pass the false belief test. In fact, all the experimental evidence seems to suggest that this is what happens to most children between 1.5 and 4 years of age.

Perhaps the best way to look at the situation is this. The Onishi and Baillargeon experiments identify an *implicit* understanding of false belief, whereas the standard false belief tasks are testing for an *explicit* understanding of false belief. By an explicit understanding I mean one that is verbally articulated and reflective, developed as part of high-level explanations of behavior in terms of beliefs and other mental states. An implicit understanding, in contrast, is predominantly practical, focused primarily on aligning one's behavior with that of others and correctly predicting how others will behave as a function of what they have or have not seen.



### **Exercise 13.5** Explain in your own words the difference between implicit and explicit understandings of false belief.

In the remainder of this chapter we will be focusing primarily on what it takes for a child to understand false belief explicitly. As we have already seen, there is evidence from (for example) pretend play suggesting that young children are capable of forms of metarepresentation considerably before they have an explicit understanding of false belief. The Onishi and Baillargeon experiments add an additional data point by showing that young children can have an implicit understanding of false belief more than 2 years earlier. This raises the very interesting question of how an implicit understanding of false belief fits into the overall development of what cognitive scientists call the mindreading system.

## 13.3 The Mindreading System

We have explored some of the connections between mindreading and pretend play. The principal link between them, according to the model first proposed by Alan Leslie and



developed by many others, is that both exploit metarepresentational skills. The model is built around the idea that mindreading and pretend play have a common information-processing structure. Both involve a “decoupling” of representations from their usual functions. In pretend play these decoupled representations serve as inputs to the PRETEND operation. In mindreading the theory of mind system uses these decoupled representations to make sense of what is going on in other people’s minds.

However, as we saw when we looked at the false belief task, some of the more complex types of mindreading emerge much later in cognitive development than pretend play, even though they both involve a sophisticated type of information processing that involves representing representations. Young children start to engage in pretend play well before they are 2 years old, but it is not until the age of around 4 that they have a rich enough understanding of belief to pass the false belief task. So, what happens in between?

## First Steps in Mindreading

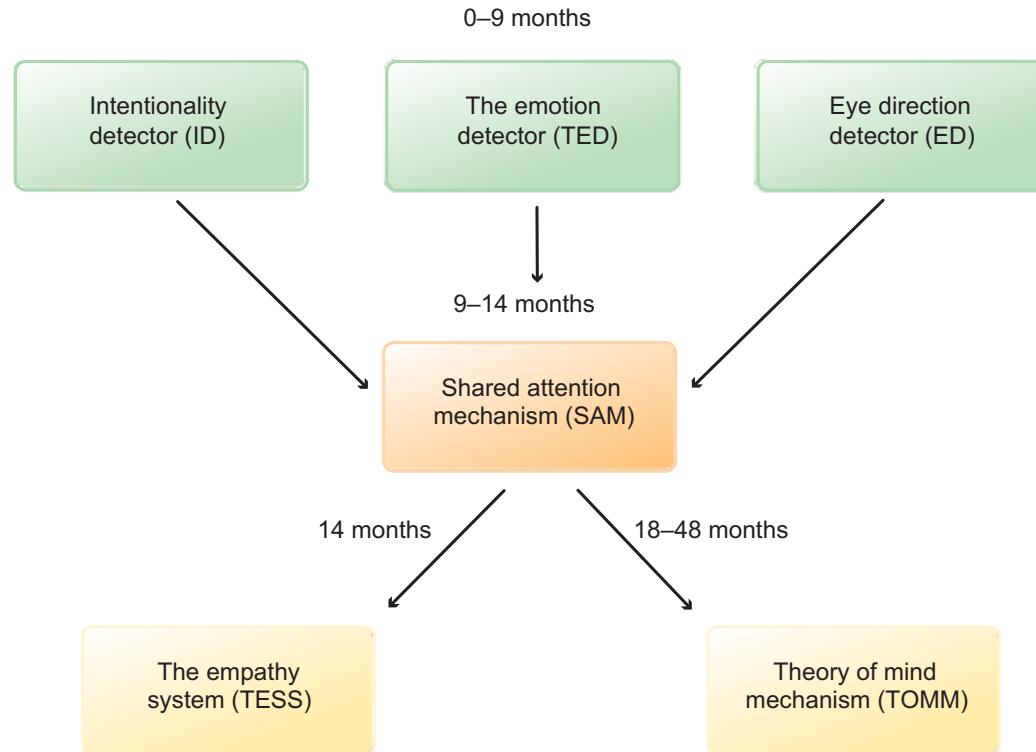
The developmental psychologist Simon Baron-Cohen was one of the co-authors of the 1985 paper that we looked at in the last section – the paper that first drew the connection between autism and problems in mindreading. Since then he has developed and fine-tuned a model of how mindreading emerges in infants and young children.

The theory of mind mechanism (TOMM) identified by Alan Leslie is the end point of the development of mindreading. But there are several stepping-stones on the way. Each of these stepping-stones opens up a different type of mindreading to the young infant. For Baron-Cohen, mindreading is a highly complex suite of abilities. It emerges in stages, with each stage building on its predecessors. The basic components of the latest version of the model are illustrated in Figure 13.6.

According to Baron-Cohen’s model, the foundations of mindreading emerge in the earliest months of infant development. The most basic mindreading skills are all in place by the time a normal infant is 9 months old. These basic mindreading skills are all essentially perceptual in nature. They involve the infant becoming perceptually sensitive to behavioral manifestations of psychological states.

The *intentionality detector* (ID) allows the infant to identify purposeful movements. When an agent makes a self-propelled movement, ID codes the movement as being goal-driven – it allows the infant to identify her mother’s arm movement as a reaching, for example. At a more fundamental level, ID allows the infant to distinguish the animate, goal-driven entities from the other objects it encounters.

A good way of finding out the apparent goal of a purposeful movement is to check where the agent is looking – since agents tend to keep their eyes on the target. So, one of the most fundamental tools for making sense of the social world is the ability to track other people’s eye movements. This is the function of the *eye direction detector* (EDD). Whereas ID enables the infant to detect purposeful movements, the job of EDD is to help the infant identify the goals of the movement. The two mechanisms are highly



**Figure 13.6** Baron-Cohen's model of the mindreading system.

complementary. There is little point in knowing that a movement is purposeful unless one has some idea what the goal is.

But there is more to making sense of people's movements than identifying purposeful movements and their goals. Young infants beginning to negotiate the social world need to be sensitive to the motivations and moods of the people with whom they are interacting – complete strangers, as well as their caregivers and other family members. This is the job of the *emotion detector* (TED). The emotion detector allows infants to understand not just that agents make movements toward particular goals, but also *why* those movements are being made and what sort of movements they are. Are they playful movements, for example, or protective ones? Sensitivity to moods and emotions is a first step toward understanding the complexities of psychology.

According to Baron-Cohen, the three basic components of the mindreading system are all in place by the time the infant is 9 months old. Well before the end of their first year, human infants can distinguish animate objects from inanimate ones. They can track where other people are looking and pick up on their moods. All of these skills are primarily perceptual. The infant is learning to pick up clues about people's psychology from what she can perceive of their physical characteristics and movements. Moods are revealed in facial expressions, and in tone of voice. Animate beings move in very different ways from inanimate objects – their movements are much less regular and much harder to predict, for example. The



orientation of the head is a good clue to eye gaze. In all these cases the infant is decoding the perceived environment in terms of some very basic psychological categories.

These three basic systems (ID, TED, and EDD) all involve relatively simple types of representation. They all involve representing other agents as having certain fairly basic features. TED, for example, involves “tagging” other agents with primitive representations of their moods (happy, sad, angry, frightened). EDD involves identifying a *dyadic* relation between an agent and an object (*Dad sees the cup*, for example). Dyadic relations have two parts. The dyadic relation of seeing is a relation between an agent and an object. ID also produces representations of dyadic relations between agents and objects. The dyadic relations here all involve intentional movements, such as reaching, or following, or pushing.

## From Dyadic to Triadic Interactions: Joint Visual Attention

Between the ages of 9 and 14 months a very important transformation takes place in the young infant’s mindreading skills. In the first 9 months of life infants are capable of understanding people and interacting with them in certain very basic ways. They are also capable of understanding objects and manipulating them. But for the very young infant these are two separate activities. Starting at the age of 9 months the infant learns to combine them. Infants become capable of employing their interactions with people in their interactions with objects, and vice versa. This is illustrated in *joint visual attention*.

Joint visual attention occurs when infants look at objects (and take pleasure in looking at objects) because they see that another person is looking at that object – and because they see that the other person sees that they are looking at the object. Joint visual attention is a collaborative activity. The infant does not just represent a dyadic relation between her mother and a cup, for example. The infant learns to represent different *triadic* (or three-way) relations between herself, the mother, and the cup – as well as to initiate them with pointing and other gestures.

In joint visual attention the infant exploits representations such as the following:

Mother SEES (I SEE the cup)  
I SEE (Mother SEES the cup)

Joint visual attention becomes possible when the infant is able to embed representations – that is, to represent that an agent (whether herself, or someone else) is representing someone else’s representation. This is very different from the type of information processing involved in detecting eye direction or sensitivity to moods. It makes possible a whole range of coordinated social behaviors in which infants and their caregivers take enormous pleasure in collaborative games – games that involve and exploit an awareness of what others are doing and how they too are participating in the game.

This distinctive kind of information processing is carried out in what Baron-Cohen has termed the *shared attention mechanism* (SAM). The emergence of the shared attention mechanism is a crucial stage in the development of the young child's mindreading skills. The connections with autism are very suggestive here too. We saw in the last section that autistic children have well-documented problems both with advanced mindreading (of the sort required for successful performance on the false belief task) and with pretend play. It turns out that autistic children also have difficulties with joint attention – and that there is a strong correlation between the severity of their social impairments and their inability to engage in joint attention.

The shared attention mechanism is also very important for language development. Pointing to objects is a very important way of teaching children what words mean. But in order for children to pick up on the cues that they are being given they need to be able to understand that they and the person pointing are jointly attending to the very same thing. Without this, children cannot understand the instructions that they are being given.

## TESS and TOMM

In Baron-Cohen's model, SAM is a crucial foundation for the final two components of the mindreading system. We have already encountered one of these components – the theory of mind mechanism (TOMM). Earlier versions of Baron-Cohen's model contained only TOMM after SAM. Recently, however, he has added an additional component, which he calls TESS (for *the empathizing system*). For normal social development it is not enough simply to be able to identify other people's emotional states and moods. The developing child needs to learn to respond appropriately to those emotional states and moods. This is where empathy comes in.

Psychosocial disorders such as psychopathy suggest that TOMM and TESS can come apart (and hence that there are two distinct and separable mechanisms carrying out the different tasks of identifying other people's mental states and developing affective responses to those mental states). Psychopaths have profound social problems, but these problems are very different from those suffered by autistic people. Psychopaths are typically very good at working out what is going on in other people's heads. The problem is that they tend not to care about what they find there – and in fact they use their understanding to manipulate other people in ways that a normal person would find unacceptable. Diagnosing psychopathy is a very complex business, but psychiatrists typically put a lot of weight on basic failures of empathy – on failure to feel sympathy when someone else is in pain or obvious distress, for example.

TESS can only emerge if the basic capacity for shared attention is in place. In many ways empathy is a matter of being able to put oneself in someone else's position – to imagine what it would be like to be someone else, and to find oneself in the situation that they find themselves in. Shared attention basically exploits the same ability, it is



just being applied in a much more limited sphere. The child engaged in joint visual attention or collaborative play is able to adopt someone else's visual perspective, to represent how things look to someone else. As they do this more and more they eventually bootstrap themselves into the ability to understand someone else's emotional perspective on the world – to understand not just how a situation looks to someone, but how that situation affects them.

The possibility of psychopathy shows (according to Baron-Cohen) that TESS and TOMM are distinct, although they both emerge from a common foundation in SAM. They develop more or less in parallel, with TESS emerging a little earlier, but TOMM taking much longer to emerge completely. The first stages in the development of TOMM are taken as early as 18 months, which is when typical young children start to engage in pretend play. But full-fledged TOMM does not emerge until much later in development – at around the age of 4, which is when young children tend on average to pass the false belief test.



## Summary

This chapter has explored the idea that there is a dedicated system for mindreading – for understanding other minds and navigating the social world. The chapter began by reviewing Leslie's theory that mindreading exploits a set of basic abilities that are also deployed in pretend play. These are abilities for metarepresentation – for representing representations. We looked at a famous set of experiments using the false belief task that seemed to show that autistic children (who are known to be deficient in pretend play) are also impaired in tasks involving reasoning about other people's beliefs. Mindreading is a complex phenomenon and we looked at a model of mindreading that sees it as made up of six distinct components, emerging at different stages in cognitive development.

## Checklist

**Alan Leslie's model of mindreading in young children is based on an analogy with the information processing involved in pretend play.**

- (1) The emergence of pretend play in the second year of life is a major milestone in cognitive and social development.
- (2) In pretend play some of an infant's *primary representations* of the world and other people become "decoupled" from their usual functions while preserving their ordinary meaning.
- (3) Leslie thinks that primary representations function in the same way in pretend play as in mindreading. Both pretend play and mindreading exploit *metarepresentation*.
- (4) Children with autism have significant problems both with mindreading and with pretend play.

- (5) The false belief task (developed by Heinz Wimmer and Joseph Perner) is a standard test of mindreading abilities in children. It tests whether children are able to abstract away from their own knowledge to understand that someone else can have different (and mistaken) beliefs about the world.

**High-level mindreading involves attributing *propositional attitudes* (such as beliefs and desires) to other people. But high-level mindreading depends upon a complex system of lower-level mechanisms – as in Simon Baron-Cohen's model of the overall mindreading system.**

- (1) The *intentionality detector* is responsible for perceptual sensitivity to purposeful movements.
- (2) The *eye direction detector* makes it easier to identify the goals of purposeful movements and to see where other people's attention is focused.
- (3) The *emotion detector* gives a basic sensitivity to emotions and moods, as revealed in facial expressions, tone of voice, etc.
- (4) The *shared attention mechanism* makes possible a range of coordinated social behaviors and collaborative activities.
- (5) The *empathizing system* is responsible for affective responses to other people's moods and emotions (as opposed to simply identifying them).

## Further Reading

Leslie first presented his metarepresentational theory of pretend play and mindreading in Leslie 1987. The theory has been considerably modified and developed since then. See Leslie and Polizzi 1998, Leslie, Friedman, and German 2004, and Leslie, German, and Polizzi 2005 for updates. For more recent research on pretend play from Leslie's group, see Friedman and Leslie 2007 and Friedman et al. 2010. For a general review of research on pretend play, see Weisberg 2015.

The false belief task discussed in the text was first presented in Wimmer and Perner 1983. It has been much discussed (and criticized). For powerful criticisms of its claimed relevance to theory of mind, see Bloom and German 2000. Perner's own theory of mindreading is presented in his book *Understanding the Representational Mind* (1993). The claim that young infants can pass a version of the false belief task was made in Onishi and Baillargeon 2005. For discussion, see Perner and Ruffman 2005, Rakoczy 2012, and Heyes 2014. Poulin-Dubois, Brooker, and Chow 2009 reviews studies on infant mindreading abilities.

Numerous recent reviews discuss both implicit and explicit false belief understanding (Baillargeon, Scott, and He 2010; Beate 2011; Low and Perner 2012; Luo and Baillargeon 2010; Perner and Roessler 2012; Trauble, Marinovic, and Pauen 2010). For a philosophical perspective on this research, see Carruthers 2013.

The idea that autism is essentially a disorder of mindreading was first presented in Baron-Cohen, Leslie, and Frith 1985. For a book-length discussion of autism as "mindblindness," see Baron-Cohen 1995. This interpretation of autism has been challenged – see, for example, Boucher



1996. The papers in Baron-Cohen, Tager-Flusberg, and Cohen 2000 discuss autism from the perspective of developmental psychology and cognitive neuroscience. For a survey, see Baron-Cohen 2009. Fletcher-Watson et al. 2014 reviews studies of clinical interventions based on the mindreading theory of autism. Roeyers and Denurie 2010 and Mathersul, McDonald, and Rushby 2013 discuss alternatives to the false belief task usable in adolescents and adults.





## CHAPTER FOURTEEN

# Mindreading: Advanced Topics

### OVERVIEW 357

<b>14.1 Why Does It Take Children So Long to Learn to Understand False Belief?</b>	358
Leslie's Answer: The Selection Processor Hypothesis	358
An Alternative Model of Theory of Mind Development	360
<b>14.2 Mindreading as Simulation</b>	363
Standard Simulationism	363
Radical Simulationism	365

<b>14.3 The Cognitive Neuroscience of Mindreading</b>	365
Neuroimaging Evidence for a Dedicated Theory of Mind System?	366
Neuroscientific Evidence for Simulation in Low-Level Mindreading?	369
Neuroscientific Evidence for Simulation in High-Level Mindreading?	373



### Overview

Mindreading is the ability to understand other people's mental states. It is key to human social interaction. Chapter 13 introduced some prominent themes in the cognitive science of mindreading. We looked at Alan Leslie's influential idea that the roots of mindreading in early childhood lie in pretend play and other activities that involve *metarepresentation* (the ability to think about thinking, as opposed to just being able to think about objects in the world). We saw how this way of thinking about mindreading is supported by the best-known test of mindreading abilities – the *false belief task*, which tests for understanding of the basic fact that people can have mistaken beliefs about the world. And then we looked at Simon Baron-Cohen's longitudinal model of mindreading, which traces how understanding of belief and other complex psychological states emerges from such more primitive abilities, such as eye gaze tracking and emotion detection.

In this chapter we continue investigating mindreading. We will tackle some more advanced topics, starting in Section 14.1 with a problem that we encountered in the last chapter. Why does it take so long for children to pass the false belief task, if (as Leslie and other believe) they are

capable of metarepresentation much earlier? We look at two different explanations – one from Leslie and one from Josef Perner (who originally developed the false belief task).

In Section 14.2 we turn to a very different way of thinking about mindreading. According to the *simulation theory*, mindreading is not really metarepresentational at all. It doesn't require dedicated information-processing systems for identifying and reasoning about other people's mental states. Instead, we make sense of their behavior by running our "ordinary" information-processing systems offline in order to simulate how other people will solve a particular problem, or react to a particular situation.

Finally, in Section 14.3 we turn to cognitive neuroscience, exploring how the techniques and technologies discussed in Chapter 9 have been used to study mindreading. In addition to studying the neural basis of mindreading (the areas in the brain that are most involved when subjects perform mindreading tasks), neuroscientists have discovered experimental evidence consistent with the simulation theory discussed in Section 14.2.

## 14.1 Why Does It Take Children So Long to Learn to Understand False Belief?

Look back at Figure 13.6 in the previous chapter. This traces the emergence of the theory of mind module (TOMM). The process is long and drawn out. It begins at around 14 months (when the infant starts to engage in pretend play) and is not complete until the child is around 4 years old (when the young child acquires the understanding of complex mental states tested in the false belief task). But why does this process take so long?

If Leslie's analysis (see Section 13.2) is correct, then information processing in the TOMM essentially exploits the machinery of metarepresentation and "decoupled" primary representations. The same machinery is involved both in pretend play and in the attribution of beliefs. So why are infants able to engage in pretend play so much earlier than they are capable of understanding beliefs and passing the false belief task?

### Leslie's Answer: The Selection Processor Hypothesis

According to Leslie and his collaborators, there is a long time lag between when the capacity for metarepresentation first emerges (during the second year) and when children generally pass the false belief test (toward the end of the fourth year), because there are actually two very different abilities here. The first is the ability to attribute true beliefs to someone else. The second is the ability to attribute false beliefs. These two abilities emerge at very different times in development. On Leslie's model, young children are able to attribute true beliefs from a relatively early age. But they are *only* able to attribute true beliefs. They can only succeed on the false belief task when they learn to "switch off," or inhibit, the default setting that other people's beliefs are true.

Leslie presents this in terms of a mechanism that he calls the *selection processor*. The selection processor is set up to favor true beliefs. So, the selection processor's default setting favors the true belief candidate. In the false belief task, this is the belief that Sally believes that the marble is in the box. But in this case, there is evidence to the contrary. The child knows



that Sally did not see the marble being moved from the basket to the box. But for this countervailing evidence to be effective, the selection processor's default setting needs to be overridden. This is what separates children who pass the false belief task from children who fail. The ones who pass are able to *inhibit* the bias in favor of the true belief candidate.

So, for Leslie, the problem is not with TOMM itself. TOMM is in place from the pretend play stage. It is just that it initially only works to attribute true beliefs. Success on the false belief task comes only when the young child acquires a more general capacity for executive control. Is there any way of testing this general hypothesis?

One way to test this hypothesis would be to alter the false belief task to make greater demands on mechanisms of inhibition and control. If the task makes greater demands on inhibitory control, and inhibitory control is the factor that explains success rather than failure on the false belief task, then one would expect that success rates on the altered task would be lower than on the original task.



### **Exercise 14.1** Explain and assess this reasoning in your own words. Can you think of other ways to test the hypothesis?

A study published by Leslie and Pamela Polizzi in 1998 reported a number of experiments adopting this general strategy. Here is a representative example. Children are presented with a scenario in which a girl (let's call her Sally, for continuity) is asked to place food in one of two boxes. The twist to the tale is that one of the boxes contains a sick kitten. Because eating the food might make the kitten worse, Sally wants to avoid putting the food into the box with the kitten in it. So, Sally has what Leslie, German, and Polizzi term an *avoidance-desire*. The point is that avoidance-desires are inhibitory. An avoidance-desire is a desire *not* to do something.

There were two conditions:

In the *true belief* condition, the kitten is moved from Box A to Box B in front of Sally.

In the *false belief* condition, the kitten is moved without Sally seeing.

Children in each condition are asked to predict which box Sally will put the food in. There is no question here about whether the children understand false belief. All the children were able to pass the standard false belief task and all of them answered correctly when they were asked where Sally thought the kitten was.

In the true belief condition, the child knows that the kitten is in Box B (since she saw the kitten being moved there) and she knows that Sally wants to avoid putting the kitten and the food in the same box. So, she needs to predict that Sally will put the food in Box A. But for that she needs to be able to make sense of Sally's inhibition of what most people would normally want to do – which is to give food to a kitten. It turned out that a very high percentage (well over 90 percent) of the children in the experiment were able successfully to predict where Sally would put the food in the true belief condition.

Now consider the false belief condition. The child still knows that the kitten is in Box B and she still knows that Sally wants to make sure that the kitten does not get the food. But now she also needs to take on board the fact that Sally did *not* see the kitten being

moved from Box A to Box B. So, as on the standard false belief task, she needs to *inhibit* her own knowledge of where the kitten is. Now the children are being asked to do two things at once – to inhibit their own knowledge of where the kitten is, as well as to make sense of Sally's inhibition of the normal desire to give food to a kitten. There is a *double inhibition* required.

According to Leslie, German, and Polizzi this is why the success rate in the false belief condition is so much lower than in the true belief condition. It turned out that only 14 percent of children in the study succeeded in the false belief condition (as opposed to 94 percent in the true belief condition). Their hypothesis is that the double inhibition places much higher demands on the selection processor than the ordinary false belief tasks.



**Exercise 14.2** Explain the reasoning behind this experiment in your own words and assess it.



**Exercise 14.3** Is the selection processor hypothesis compatible with the Onishi and Baillargeon data suggesting an implicit understanding of false belief in 15-month-old infants? If so, how? If not, why not?

## An Alternative Model of Theory of Mind Development

Here is an alternative way of thinking about the problem, due to the developmental psychologist Joseph Perner (one of the two authors of the original paper that presented the false belief task).

Perner's thinking is very much informed by influential theories in philosophy about the nature of belief, and other mental states that philosophers collectively label *propositional attitudes*. Belief is called a propositional attitude because it involves a thinker taking an attitude (the attitude of belief) toward a proposition. So, if I believe that it is now raining, then I am taking the attitude of belief to the proposition *it is now raining*.

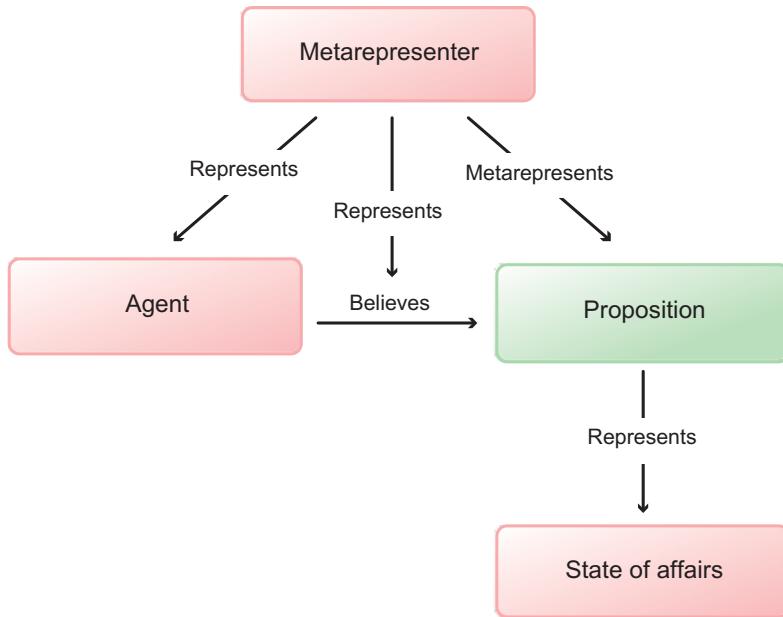
What are propositions? We can simply think of propositions as representations of the world that can be either true or false. This means that if a young child (or anyone else, for that matter) is to attribute a belief to someone else, she must represent that person as standing in the belief relation to a representation of the world *that can be either true or false*.

Understanding propositions in this way, Perner rejects Leslie's claim that there could be a theory of mind mechanism that only attributes true beliefs. Leslie may well be right that young children are attributing to others some sort of psychological state that is always true (from the child's perspective). But, according to Perner, that psychological state cannot be the state of belief. Beliefs are just not the sort of thing that can always be true.



**Exercise 14.4** State in your own words and assess Perner's objection to Leslie's model of the TOMM.

In fact, Perner draws a stronger conclusion. If the psychological states that the child attributes are always true (from the child's perspective), then the child is not really engaged in metarepresentation at all. The child is certainly representing another person as being in



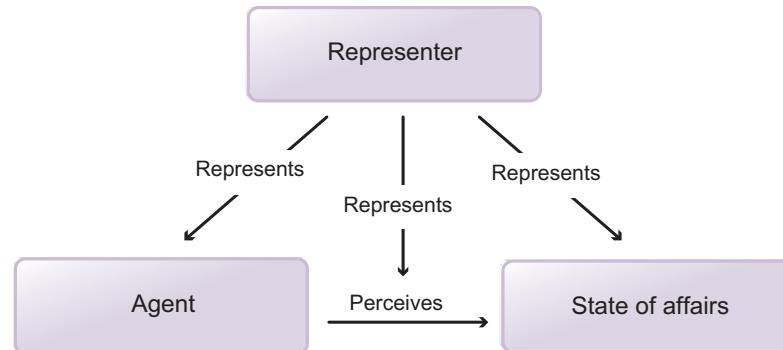
**Figure 14.1** What goes on when one subject represents another's belief. Note that representing belief requires metarepresentation.

a psychological state. But they can do that without engaging in metarepresentation. Since the content of the psychological state tracks what the child considers to be the state of the world, the child does not need to deploy any resources over and above the resources that she herself uses to make sense of the world directly.

Compare belief with perception. As we saw when we first encountered the false belief task in Section 13.2, perception is *factive*. I can only perceive that things are a certain way if they really are that way. I can only perceive that it is raining if it really is raining. The factive nature of perception carries across to what happens when we represent someone else as perceiving something. I cannot represent someone else as perceiving that it is raining unless I think that it really is raining. But this does not hold for belief.

Figures 14.1 and 14.2 make it easier to see what is going on here. Figure 14.1 shows full-blown metarepresentation. The structure of metarepresentation is triadic. A metarepresenting subject has to represent another psychological subject, a psychological relation such as belief, and the proposition that is believed. The proposition represents a particular state of affairs in the world (the state of affairs of the marble being in the box, or of the cat being on the mat). It can be either true or false. But, as far as the accuracy of metarepresentation is concerned, what matters is not whether or not the proposition is true, but rather whether or not the metarepresenting subject has identified it correctly.

In Figure 14.2, in contrast, we see what is going on when a psychological subject represents another psychological subject's perceptual state. Here there is no need for a proposition to be identified. There is no need for metarepresentation in the strict sense. All that the person representing the psychological state needs to do is to represent directly a relation between the perceiver and the state of affairs in the world that the perceiver is perceiving.



**Figure 14.2** What goes on when one subject represents another's perception. Note that representing perception does not require metarepresentation.

On Perner's view of mindreading, therefore, metarepresentation in the full sense of the word does not appear until fairly late in development. It is only when children start to understand the possibility of false belief that we see the emergence of what Perner calls the *representational mind*.

But now we have a new problem. If Perner is right, however, that metarepresentation does not emerge until children pass the false belief test, then we need to find another way of interpreting what is going on in pretend play.

Perner's key idea here is that primary representations can be "decoupled" from reality without there being metarepresentation going on. Thinkers can decouple primary representations from reality without representing them as representations.

*Counterfactual thinking* is a good example. We engage in counterfactual thinking when we think about how things might be (but are not). If I am wondering whether I have made the right choice of restaurant I might start to think about how things might have turned out had I made a different choice. I might imagine how things would be in a different restaurant, for example – the different things that I could have ordered, the different clientele, the lack of noise, and so on. The representations in counterfactual thinking are decoupled in the sense that they are being used to think about how things might be, rather than about how they are. But they do not involve metarepresentation. When I think about the steak that I might now be having in the restaurant over the street, my representation of the steak is decoupled from reality (because, after all, I am not thinking about any existing steak). But I am not engaged in metarepresentation – I am thinking about the steak that I could be having, not about my representation of the steak.

Metarepresentation is a matter of *thinking about* decoupled representations (thinking that is directly focused on representations, rather than on the world). But counterfactual thinking is a matter of *thinking with* decoupled representations (using decoupled representations to think about the world). When the child pretends that the banana is a telephone, she is decoupling her primary representations of the telephone and applying them to the banana. But at no point is she representing those primary representations – and so she is not engaged in metarepresentation.



A cognitive scientist who adopts Perner's interpretation of pretend play can nonetheless adopt many of Leslie's specific proposals about the information processing in pretend play. She could also adopt the model of the complete mindreading system proposed by Simon Baron-Cohen (although the emergence of the TOMM would have to be dated somewhat later). Because of this, one might well think that there is much more agreement than disagreement between Leslie and Perner. In fact, this turns out to be exactly right when we look at a very different model of mindreading that some cognitive scientists and developmental psychologists have proposed. This is the *simulationist* model that we will examine in the next section. From the perspective of simulationists, Leslie and Perner are really both thinking about mindreading in the same (mistaken) way.



**Exercise 14.5** Go back to Section 13.1 and identify how Leslie's basic model of pretend play would need to be modified in order to accommodate Perner's interpretation.



**Exercise 14.6** Is Perner's interpretation compatible with the Onishi and Baillargeon data suggesting an implicit understanding of false belief in 15-month-old infants? If so, how? If not, why not?

## 14.2 Mindreading as Simulation

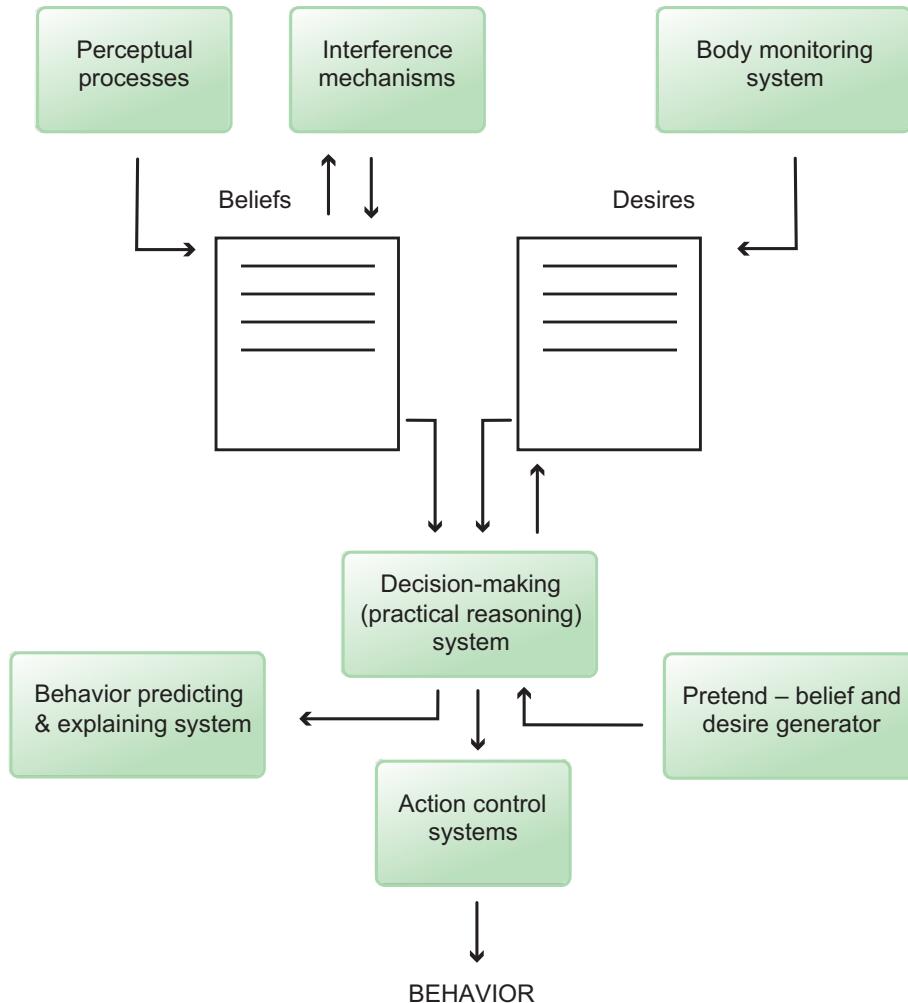
The last section focused primarily on the differences between the models of mindreading developed by Alan Leslie and Joseph Perner. But still, they have a lot in common. They both think that mindreading is basically a theoretical accomplishment. It requires bringing a specialized body of knowledge (theory of mind) to bear in order to explain and predict the behavior of others.

This section explores an alternative to this theory-based approach. According to the *simulation theory*, there is no specialized theory of mind mechanism. Instead, theory of mind processing is carried out by the very same systems that are responsible for ordinary decision-making and for finding out about the world.

Here is the basic idea. Suppose that we have a reasonable sense of the beliefs and desires that it would be appropriate to attribute to someone else in a particular situation. To find out how they will behave, we don't use specialized knowledge about how mental states typically feed into behavior. Instead, we use our own decision-making processes to run a simulation of what would happen if we ourselves had those beliefs and desires. We use the output from that simulation process to predict how the person will behave.

### Standard Simulationism

One way of developing this basic idea was originally proposed by the developmental psychologist Paul Harris and subsequently developed by the philosopher Alvin Goldman. We can call their theory *standard simulationism*. According to standard simulationism, the process of simulation has to start with the mindreader *explicitly* (although not necessarily consciously) forming judgments about how the other person represents the relevant situation and what



**Figure 14.3** A schematic version of standard simulationism. Note that the ordinary decision-making system is being run offline with pretend inputs. (Adapted from Nichols et al. 1996)

they want to achieve in that situation. These judgments serve as the input to the ordinary decision-making system. This general model is illustrated in Figure 14.3.

Where do these judgments come from? Here's Goldman's proposal. He thinks that we work by analogy with our own beliefs and desires. We know which beliefs we tend to form in response to particular situations. And so, we assume that others will form the same beliefs, unless we have specific evidence to the contrary – e.g., about how they have acted in the past, or about temperamental differences, or the different information that they might have. When we do have such additional evidence, we make adjustments by thinking about what we would do if we had those temperamental features or that extra information.

So, on Goldman's model, knowledge of others rests upon self-knowledge. And he thinks that we have a special mechanism for finding out about our own beliefs, desires, and other propositional attitudes – a self-monitoring mechanism that philosophers call *introspection* or *inner sense*.



## Radical Simulationism

There is a second way of developing the basic simulationist idea – often called *radical simulationism*. Radical simulationism has been developed primarily by the philosophers Robert Gordon and Jane Heal. The intuitive idea behind radical simulationism is that, instead of coming explicitly to the view that the person whose behavior I am trying to predict has a certain belief (say, the belief that *p*), what I need to do is to imagine how the world would appear from her point of view.

Here is the difference. In standard simulationism, the process of simulation starts with beliefs about another person's beliefs and desires. They are metarepresentations. According to radical simulationism, on the other hand, what the simulator is thinking about is the world, rather than the person they are simulating. The simulator is thinking about the world *from the perspective of the person being simulated*, rather than thinking about their beliefs, desires, and other psychological states.

Radical simulationism proposes mindreading without metarepresentation. This is because it is world-directed, rather than mind-directed. And as a result, it gives a very different account of what is going wrong when children fail the false belief test. For the radical simulationist, children who fail the false belief test lack imaginative capacities. Their capacity to project themselves imaginatively into someone else's position is not sufficiently developed. They are not yet able to form beliefs *from a perspective other than their own*. They are capable of imaginative perceiving. That is, they can adopt someone else's perceptual perspective on the world – they can appreciate how things *look* to Sally. But they are not capable of imaginatively working their way into the beliefs that someone might have about the world.



**Exercise 14.7** We have now looked at four different ways of thinking about the false belief task. Draw up a table indicating the four different proposals that have been made for explaining what it is that the false belief task is testing for.

## 14.3

### The Cognitive Neuroscience of Mindreading

This section explores what we can learn about mindreading from cognitive neuroscience. The neuroscience of mindreading has become a very hot topic in recent years and we can only scratch the surface here. But we can focus the issues by concentrating on three questions that emerge from our earlier discussion.

Several of the models of mindreading that we have been looking at hold that the mind contains a dedicated theory of mind mechanism (TOMM). This is the information-processing system responsible for reasoning about other people's beliefs, desires, and other propositional attitudes. So, a natural question to ask is,

*Question 1* Is there any evidence at the neural level for the existence of a TOMM?

We also looked at the alternative proposal that mindreading is an exercise in simulation. It does not exploit systems specialized for mindreading. Instead, the information

processing in mindreading is carried out by cognitive systems that also do other things. We can call these *co-opted systems*. So, we can ask whether there is any evidence at the neural level for this way of thinking about mindreading.

There are really two different questions here. We can make a distinction between low-level mindreading and high-level mindreading. Low-level mindreading involves detecting emotions, identifying goal-driven actions, sensitivity to eye gaze, and so on. High-level mindreading involves identifying and reasoning about beliefs, desires, and other psychological states. This gives us two further questions:

*Question 2* Is there evidence at the neural level that *low-level mindreading* is a process of simulation involving co-opted systems?

*Question 3* Is there evidence at the neural level that *high-level mindreading* is a process of simulation involving co-opted systems?

## Neuroimaging Evidence for a Dedicated Theory of Mind System?

In Chapter 11 we saw that neuroimaging allows cognitive scientists to map activity in the brain while subjects are performing specific tasks. So, to use neuroimaging to investigate whether there is a dedicated TOMM, experimenters have looked for brain regions with these two characteristics:

- 1 They show increased activity in response to information-processing tasks that require the subject to attribute beliefs.
- 2 These increased activation levels are specific to tasks involving belief attribution – as opposed, for example, to reflecting demands on general reasoning, or the fact that people (rather than inanimate objects) are involved.

As far as (1) is concerned, it is very important that a candidate TOMM region should show increased activation both for false belief tasks and for true belief tasks. What (2) is asking for is evidence that the neural systems are engaged in domain-specific processing. In order to establish that (2) holds, experimenters need to make sure that they have controlled for domain-general processes (such as language or working memory).

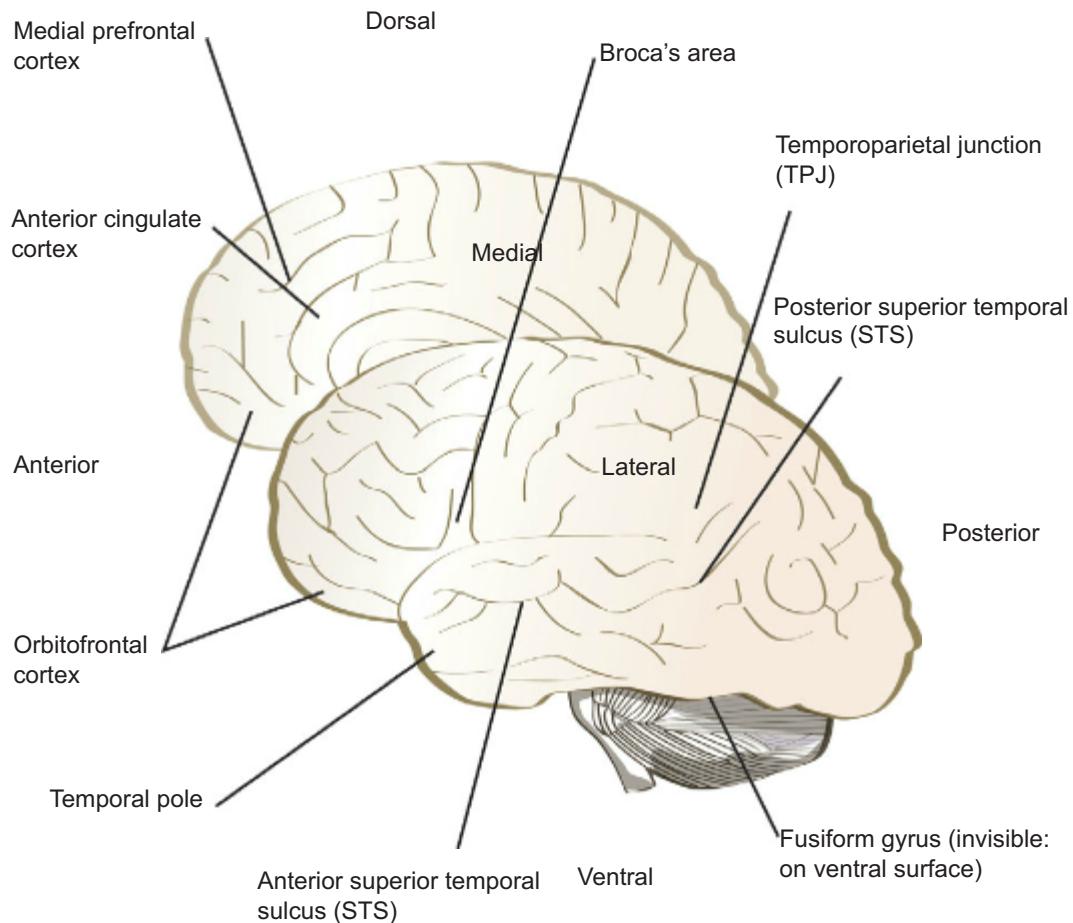
Neuroimaging studies have identified a number of brain regions as showing increased activation in tasks that seem to require reasoning about beliefs. Most of these studies have involved versions of the false belief test, although some have explored different paradigms. The cognitive psychologist Vinod Goel, for example, ran a series of studies in which he asked subjects to decide whether Christopher Columbus would have been able to work out the function of an object from a picture – the idea being that this task requires subjects to reason about the sort of beliefs that a fifteenth-century explorer would have been likely to have. Other studies had subjects read a short story and then answer questions on it. Some of the questions required making inferences about the beliefs of characters in the story and others not.



Studies such as these have identified a number of brain regions as potentially forming part of a dedicated theory of mind system. These include (working more or less from front to back):

- medial prefrontal cortex
- anterior cingulate cortex
- orbitofrontal cortex
- temporal pole
- Broca's area
- anterior superior temporal sulcus
- fusiform gyrus
- temporoparietal junction
- posterior superior temporal sulcus

This is a long list and, as we see in Figure 14.4, these regions collectively cover a large area of the brain.



**Figure 14.4** Schematic representation of brain regions associated with the attribution of mental states. (Adapted from Saxe, Carey, and Kanwisher 2004)

The list includes a number of brain areas thought to be specialized for other information-processing functions. Broca's area, for example, is widely held to be involved in aspects of language processing, while the fusiform gyrus includes the fusiform face area (which has been hypothesized as a dedicated face-processing system). This is not particularly surprising. The various tasks that have been used to explore belief attribution inevitably bring other capacities and abilities into play. In order to narrow the list down we need to see which (if any) of these areas satisfy (1) and (2) above.

The first stage, corresponding to (1), is to check whether particular neural regions show activation both in false belief and in true belief conditions. This is particularly important, since many neuroimaging studies follow the developmental studies in focusing only on false belief conditions. This emphasis on false belief is fine for looking at the development of mindreading in children – since the crucial developmental measure is standardly taken to be success on false belief tasks. But if we are looking for a neural substrate for belief reasoning we need to consider true belief conditions as well as false ones – after all, perhaps some of the activation in the false belief condition is due to the falsity of the belief attributed, rather than to its being a belief.

Rebecca Saxe and Nancy Kanwisher carried out a set of false belief experiments with a true belief condition as a control. We will look at these experiments in more detail below (in the context of identifying mechanisms specialized for theory of mind tasks). For the moment we need only note what happened when they did a more detailed statistical analysis of the patterns of activation within individual subjects. They found three brain regions where both true and false belief attribution tasks elicited activation in the very same voxels. (Recall that a voxel is a *volumetric pixel* representing a small volume within the brain.) These regions are the

- medial prefrontal cortex (MPFC)
- superior temporal sulcus (STS)
- temporoparietal junction (TPJ)

Applying (1), then, significantly narrows down the field. What happens when we apply (2)?

To apply (2) we need to control for other types of domain-general information processing that might be generating activation in the candidate areas. Saxe and Kanwisher introduced two control conditions, based on their analysis of what is required in order to succeed on tasks involving belief attribution.

Their first observation is that when we attribute beliefs to other people we are effectively identifying hidden causes. This is because we typically attribute beliefs when we are trying to explain or predict behavior, and we cannot do so in terms of what is immediately observable. So, in order to make sure that activation in the candidate theory of mind areas really does reflect domain-specific theory of mind reasoning, we need to rule out the possibility that what is going on is really just domain-general reasoning about hidden causes.

To do this, Saxe and Kanwisher developed a set of stories depending on nonpsychological hidden causes. Here are two:



- The beautiful ice sculpture received first prize in the contest. It was very intricate. Unfortunately, the temperatures that night hit a record high for January. By dawn, there was no sculpture.
- The night was warm and dry. There had not been a cloud anywhere for days. The moisture was certainly not from rain. And yet, in the early morning, the long grasses were dripping with cool water.

Call this the hidden causes condition.

Saxe and Kanwisher also wanted to rule out the possibility that activation is due to general reasoning about false representations – as opposed to false *beliefs*. There is nothing psychological about a false representation such as a misleading map, for example. In order to rule out the possibility that the neural areas active in belief attribution are specialized for information processing to do with representations in general rather than theory of mind, Saxe and Kanwisher used a version of the *false photograph task* originally proposed by the developmental psychologist Debbie Zaitchik.

Here is a false photograph version of the false belief task. As before, the subject is presented with a story in which Sally places a marble in the basket. A photograph is taken of the contents of the basket and placed face down. After the photograph is taken, Anne moves the marble from the basket to the box. Everything is exactly as in the false belief task – except that the subjects are asked where the object appears in the photograph. The idea behind the task is that a subject who does not understand the possibility of false representations will think that the object's location in the photograph will be where it really is – and so the photograph will depict the marble as being in the box.



### Exercise 14.8 Assess the reasoning behind the false photograph task.

Experimental subjects were presented with a number of short stories and questions in each of the three conditions. Saxe and Kanwisher found that there was significant activation in the three regions identified earlier (MPFC, STS, and TPJ) in the belief attribution condition, but not in the false representation or hidden causes conditions.

Still, all this shows is that there are brain regions specialized for processing information about mental states such as belief. These experiments do not tell us much, if anything, about what is going in those regions.

## Neuroscientific Evidence for Simulation in Low-Level Mindreading?

Look back at Simon Baron-Cohen's model of the mindreading system in Figure 13.6. The theory of mind mechanism (TOMM) is a relatively small part of the overall mindreading system – just one out of six components. At least until recently, this part of the mindreading system has received by far the most attention from cognitive scientists. This is not very surprising, since it is in many ways the most sophisticated and visible tool that we have for navigating the social world. But, as the model brings out, we have a range of other

tools besides explicit reasoning about beliefs, desires, and other propositional attitudes. So, for example, we are sensitive to other people's emotional states, to where their eyes are directed, and to what the targets of their actions are.

Simulation theorists claim that mindreading is carried out by what they call *co-opted mechanisms*. These are information-processing systems that normally serve another function and that are then recruited to help make sense of the social world. A number of experiments have been interpreted by simulation theorists as showing that co-opted mechanisms play a fundamental role in mindreading.

One very basic form of mindreading is the ability to read emotions off perceptible bodily states. Facial expressions are the most obvious example, but tone and timbre of voice are often good guides to emotions, as are global features of posture (the vitality and energy of someone's movements, for example). Young children start to develop their skills in this form of mindreading at a very early age. It is an automatic and unconscious process for normal people – fundamental to our interactions with other people, and of course to how we respond to pictures and films. In Baron-Cohen's model it is the job of a dedicated component: the emotion detector.

On the simulationist view, the emotion detector is likely to be a co-opted mechanism (or set of mechanisms). What sort of co-opted mechanisms? The most obvious candidates are the very same mechanisms that allow people to experience emotions. The simulationist approach to mindreading holds that there is a single set of emotion mechanisms that come into play both when agents are experiencing emotional states and when they detect emotions in others. Is there any evidence that this is so? Some suggestive results have come from the study of brain-damaged patients.

Here is an example. Many studies have found that a region of the temporal lobe known as the *amygdala* plays an important role in mediating fear. The experience of disgust, in contrast, is much more closely correlated with activity in the *insula*, which lies in the lateral sulcus, separating the temporal lobe from the parietal lobe. (Both the amygdala and the insula form part of the *limbic system*.)

According to simulation theory, the very same mechanism that mediates the experience of a particular emotion is recruited when the subject recognizes that emotion in someone else. So, for example, a simulationist would expect the amygdala to be active both when someone is undergoing fear and when they identify that fear in others. And, conversely, a simulationist would expect damage to the amygdala to result in a patient having problems *both* with the experience of fear and with identifying fear in others. The prediction, therefore, is that damage to brain regions that play a significant role in mediating particular emotions will result in *paired deficits* – in problems with experiencing the relevant emotion and in identifying it in others.

There is evidence of paired deficits for several different emotional states:

- **Fear:** Ralph Adolphs and his colleague have studied a number of patients with damage to the amygdala. The patient S.M., for example, had her amygdala destroyed on both sides of the brain by Urbach–Wiethe disease. She is, quite literally, fearless – although she knows what fear is, she does not experience it. She is also significantly impaired on tests that



require identifying fear on the basis of facial expression. Psychopathic patients are known to have both smaller amygdalas than normal subjects and reduced capacities for experiencing fear. It turns out that they are also much less good than normal controls at identifying fear in others.

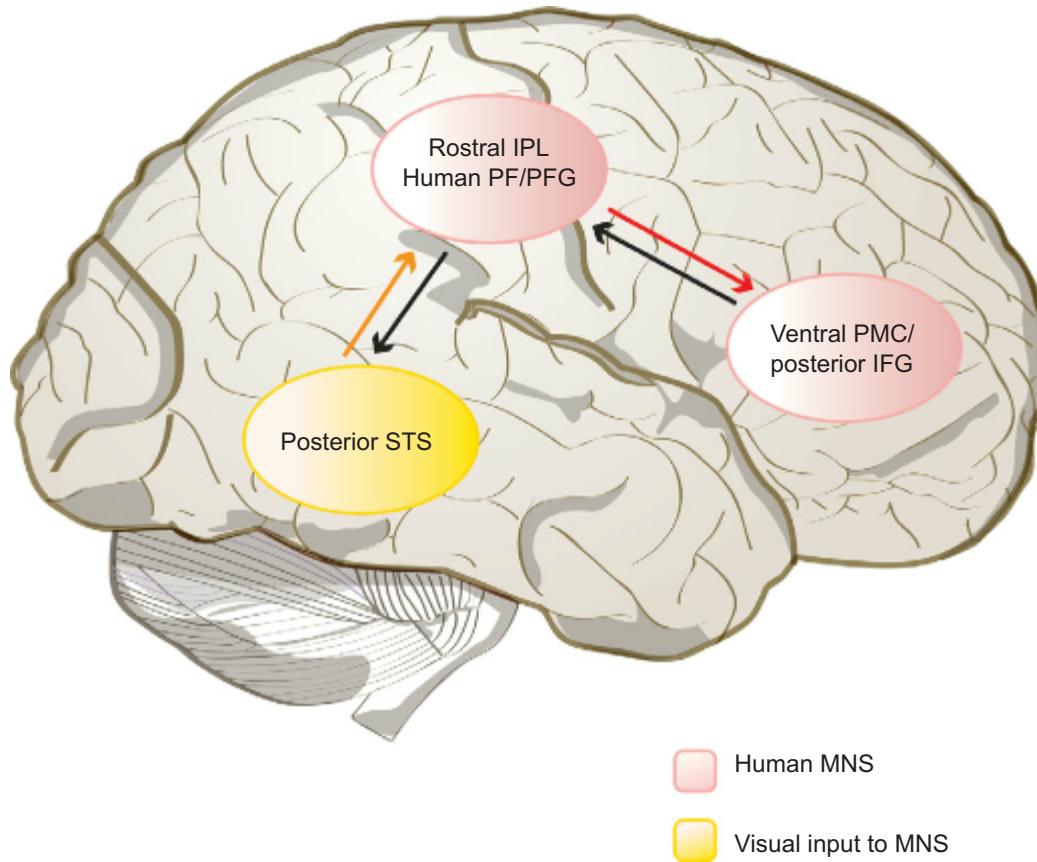
- **Anger:** The neurotransmitter *dopamine* is thought to play an important role in the experience of anger. Experiments on rats, for example, have shown that levels of aggression can be directly manipulated by raising/lowering the rat's dopamine levels. In humans, dopamine production can be temporarily blocked with a drug called *sulpiride*. Experiments have shown that subjects whose dopamine levels have been lowered in this way are significantly worse than controls in recognizing anger from facial expression – but do not have problems with other emotions.
- **Disgust:** The brain area most associated with the experience of disgust is the insula. Neuroimaging studies have shown that this area is also activated when subjects observe facial expressions of disgust. This result is confirmed by studies of brain-damaged patients. N.K., a much-studied patient suffering from damage to the insula and basal ganglia, has severe problems both in experiencing disgust and in recognizing it in others. He performs no differently from controls, however, with regard to other basic emotions (such as surprise and fear).

Supporters of the simulationist approach to mindreading have also found evidence for co-opted mechanisms in some much-publicized experiments on “mirror neurons.” We looked briefly at mirror neurons in Section 9.2, as an example of what we can learn from recording electrical activity in single neurons. (This would be a good moment to look back at Figure 9.5 to see mirror neurons in action.)

Mirror neurons were first discovered in macaque monkeys by an Italian research group led by Giacomo Rizzolatti in the mid-1990s. Rizzolatti and his colleagues were recording the responses of neurons that showed selective activation when the monkey made certain hand movements (such as reaching for a piece of food) when they noticed completely by chance that the same neurons fired when the monkey saw an experimenter making the same movement.

In monkeys the mirror neuron system is located in area F5 in the ventral premotor cortex, as well as in the inferior parietal lobe. There has been considerable discussion about whether mirror neurons exist in humans. No mirror neurons have ever been directly detected in humans – not surprisingly, since it is not usually possible to make single-cell recordings in humans. The evidence for mirror neurons in humans comes primarily from fMRI studies. Studies have found a brain system that appears to have the basic “mirroring” feature – that is, its elements show activation both when the subject performs certain actions and when others are observed making that action. Researchers have dubbed this system the mirror neuron system. The mirror neuron system is illustrated in Figure 14.5 and described in the accompanying caption.

Some cognitive scientists have suggested that the mirror neuron system functions as an *empathy* system. It allows people to *resonate* to the psychological states of other people. So, for example, studies have shown that areas in the mirror neuron system are activated both



**Figure 14.5** Schematic overview of the frontoparietal mirror neuron system (MNS) (pink) and its main visual input (yellow) in the human brain. An anterior area with mirror neuron properties is located in the inferior frontal cortex, encompassing the posterior inferior frontal gyrus (IFG) and adjacent ventral premotor cortex (PMC). A posterior area with mirror neuron properties is located in the rostral part of the inferior parietal lobule (IPL) and can be considered the human homolog of area PF/PFG in the macaque. The main visual input to the MNS originates from the posterior sector of the superior temporal sulcus (STS). Together, these three areas form a “core circuit” for imitation. The visual input from the STS to the MNS is represented by an orange arrow. The red arrow represents the information flow from the parietal MNS, which is mostly concerned with the motoric description of the action, to the frontal MNS, which is more concerned with the goal of the action. The black arrows represent reference copies of motor predictions of imitative motor plans and the visual description of the observed action. (Adapted from Iacoboni and Dapretto 2006)

when the subjects feel pain and when they observe a loved one undergoing a painful stimulus. In terms of the models that we have been using, this would mean that the mirror neuron system could serve as a neural substrate both for TED (the emotion detector system) and TESS (the empathy system). And, as the caption to Figure 14.5 brings out, it is also thought that the mirror neuron system is part of what makes imitation possible.

Some of the stronger claims that have been made in this area should be treated with caution. Quite apart from any skepticism about whether there actually are any mirror



neurons in humans, there are definite limits to the explanatory power of mirror neurons. Macaque monkeys are not very sophisticated mindreaders, to put it mildly, and so one might reasonably wonder about the role that can be played in mindreading by neural mechanisms present both in humans and monkeys.

The most likely application for mirror neurons is the information processing associated with understanding basic forms of goal-driven action – what Baron-Cohen calls the intentionality detector. There is some evidence that mirror neurons are sensitive to goals (rather than simply to bodily movements). A study published in 2001 by Alessandra Umiltà and colleagues showed that mirror neurons fire even when the monkey cannot see the final stages of the action. They used a screen to hide the experimenter's hand when it actually grasped the object and found that about 50 percent of the mirror neurons usually tuned to grasping actions were activated even in the absence of the usual visual cues for grasping. It seems that mirror neurons are sensitive to fairly abstract properties of movements – to the fact that they are goal-directed, rather than simply to their physical and observable characteristics.

In any event, mirror neurons in monkeys are direct examples at the most basic neural level of mechanisms that show the *dual purpose* structure at the heart of the simulationist approach to mindreading. And much of the evidence that has been produced in support of the existence of a mirror neuron system points to the existence of brain regions that serve both first-person and third-person roles. They are active both when the subject performs certain actions and/or undergoes experiences of a certain type – and when others are observed performing those actions and/or undergoing those experiences.

## Neuroscientific Evidence for Simulation in High-Level Mindreading?

There is far less direct evidence for simulation in high-level mindreading than in the lower-level processes that we have just been discussing. Nonetheless, there are some suggestive results.

Simulationists differ on how exactly the process of simulation is supposed to work. For standard simulationists, the process of simulation requires some form of *inference by analogy*. In essence, the simulator works out what she would do in a given situation and then infers (analogically) that the person she is trying to predict will do the same thing. Radical simulationists, in contrast, think that simulation can take place without this type of inference from oneself to others. They hold that simulation is fundamentally a matter of adopting another person's perspective – putting oneself into their shoes, as it were.

There is a prediction here. If standard simulation is a correct way of thinking about mindreading then mindreading should be both a first-person and a third-person process. The basic engine of simulation is the simulator running her own decision-making processes offline and identifying her own mental states. The results of this first-person simulation are then applied to the person being simulated. The prediction from standard simulation, therefore, is that regions of the brain specialized for what is sometimes called

*self-reflection* (i.e., identifying one's own psychological attributes, abilities, and character traits) will be active during tasks that require mindreading.

There is some evidence bearing this prediction out. A number of studies have shown that self-reflection tasks elicit activation in an area of the brain thought to be involved in high-level mindreading – the medial prefrontal cortex (MPFC – illustrated in Figure 14.4). So, for example, in one set of studies (published by William Kelly and collaborators in 2002) subjects were presented with various written adjectives and asked some questions about them. These questions were either perceptual (“Is this adjective written in italics?”), self-directed (“Does this adjective describe you?”), or other-directed (“Does this adjective describe the President?”). The self-directed questions consistently generated greater activation in MPFC.

Further support for this apparent connection between self-reflection and mindreading came from a study published by Jason Mitchell, Mazharin Banaji, and Neil Macrae in 2005. The experimenters scanned subjects while they were presented with photographs of other people and asked questions about them. Some questions required mindreading (“How pleased is this person to have their photograph taken?”), while others did not (“How symmetrical is this person’s face?”). After a short delay the subjects were presented with the photographs again and asked how similar they thought the other person was to themselves. This question is important for simulation theorists because simulation is likely to work best for people whom one thinks are similar to oneself.

These experiments produced two significant results. First, they provided further evidence that MPFC is important in high-level mindreading – MPFC showed much higher activation levels on the mindreading version of the task than on the other version. More significant was what happened when the experimenters compared activation in MPFC on the mindreading version of the task with the subjects’ subsequent judgments when they were asked how similar they perceived the other person to be to themselves. It turned out that there was a significant correlation between activation in MPFC while the subjects were answering the mindreading questions and the degree of similarity that subjects subsequently observed between themselves and the person in the photograph. The greater the perceived similarity with the person in the photograph, the higher the level of activation in the subject’s MPFC.

The cognitive neuroscience of mindreading is clearly a fascinating and thriving area. We have reviewed a number of important findings and experiments. It is far too early to draw any definite conclusions. But even this brief review illustrates very clearly two important points that emerged in earlier chapters:

- The cognitive neuroscience of mindreading involves careful calibration of results from different technologies. This comes across very clearly in the way experimenters have worked through the potential implications of mirror neurons for thinking about mindreading in humans. Single-neuron studies in monkeys have been calibrated by functional neuroimaging in humans.
- Neuroscientists interested in mindreading are not simply exploring the neural implementation of cognitive information-processing models developed in abstraction



from details about how the brain works. It is true that much of the discussion in this area is driven by psychological experiments such as the false belief task and the cognitive models that have been produced in response to them, but participants at all levels in the debate clearly recognize that techniques from neuroscience have a crucial role to play in testing, confirming, and developing cognitive models.



## Summary

This chapter explored advanced topics in mindreading. We started out with a puzzle about the emergence of mindreading. Children typically do not pass the false belief task until around 4 years of age, but other parts of the mindreading system are in place well before that. Why the delay? We looked at two possible explanations – Leslie's *selection processor* hypothesis and Perner's account in terms of metarepresentation.

Both Perner and Leslie think that there is a theory of mind module dedicated to identifying and reasoning about other people's mental states, making it possible to understand why they are doing what they are doing, and to predict what they will do in the future. We then explored an alternative model of mindreading. According to simulationists, there is no dedicated mindreading system. Instead mindreading is carried out by our "ordinary" cognitive systems running offline with pretend inputs. We looked at two different versions of simulationism – standard simulationism and radical simulationism.

Finally, we reviewed a range of evidence from cognitive neuroscience. Some neuroimaging experiments reveal areas with increased activation during mindreading tasks. At the same time, there is evidence for the simulationist hypothesis that mindreading involves co-opting ordinary decision-making systems, both from research on mirror neurons in monkeys and also from imaging studies showing that areas specialized for self-reflection are active in mindreading tasks.

## Checklist

**Young children do not typically pass the false belief task before the age of 4, although other parts of the mindreading system come onstream much sooner. Different explanations have been given of this time lag.**

- (1) Leslie argues that the theory of mind mechanism emerges during the infant's second year. But its default setting is to attribute true beliefs. Overcoming that default setting requires the emergence of an inhibitory mechanism that he calls the *selection processor*.
- (2) Support for the selection processor interpretation comes from double *inhibition experiments*.
- (3) For Perner, in contrast, children do not understand belief, properly speaking, until they pass the false belief task. Understanding belief requires the possibility of metarepresentation, and an inability to metarepresent explains failure on the task.
- (4) Perner (and others) have developed accounts of pretend play on which it does not involve metarepresentation.

Perner and Leslie (and many other cognitive scientists) are committed to the idea that there is a dedicated theory of mind system responsible for identifying and reasoning about other people's beliefs, desires, and other propositional attitudes. This basic assumption is challenged by the simulationist approach to mindreading.

- (1) Simulationists think that mindreading is carried out by "ordinary" information-processing systems that are co-opted for mindreading. We use our own mind as a model of someone else's mind.
- (2) According to standard simulationism, we predict other people's behavior, for example, by running our decision-making processes offline, with pretend beliefs and desires as inputs.
- (3) Radical simulationists hold that mindreading does not involve representing another person's psychological states. Rather, it involves representing the world from their perspective.

Cognitive neuroscientists have used a range of techniques, including single-neuron recording and functional neuroimaging, in order to test and refine cognitive models of mindreading. These are early days in the cognitive neuroscience of mindreading, but some suggestive results have already emerged. For example:

- (1) Neuroimaging studies have identified a number of brain areas that show increased activation during mindreading tasks. Experiments by Saxe and Kanwisher, for example, have highlighted the medial prefrontal cortex, the superior temporal sulcus, and the inferior parietal lobule. This is consistent with the claim that there is a dedicated theory of mind system.
- (2) There is evidence that co-opted mechanisms are used in low-level mindreading (as predicted by the simulation theory). Areas active during the experience of basic emotions such as fear, disgust, and anger are also active when those emotions are identified in others.
- (3) Mirror neurons in area F5 of the macaque brain respond both when the monkey performs an action and when the monkey observes an experimenter or conspecific perform that action. A number of researchers have hypothesized a mirror neuron system in the human brain. This may play an important role in understanding goal-directed action.
- (4) There is evidence consistent with the simulation-driven processing in high-level mindreading. Experiments have shown that areas specialized for self-reflection are also implicated in mindreading (as predicted by standard simulationism).

## Further Reading

The selection processor hypothesis was proposed in Leslie and Polizzi 1998 and then further developed in Leslie, German, and Polizzi 2005. The hypothesis is criticized in Doherty 1999. For Perner's view, see Perner, Leekam, and Wimmer 1987, Perner 1993, and Perner and Roessler 2012. Both Leslie and Perner are predicated on the assumption that children do not typically pass the false belief task before they are 4 years old or so. For a different view, claiming that success appears much earlier in development (even on the traditional version of the task), see Setoh, Scott, and Baillargeon 2016 and (for a review) Scott and Baillargeon 2017.

Mindreading was one of the earliest fields to see sustained interactions and collaborations between philosophers and psychologists. A number of influential early papers, including Heal 1986 and Gordon 1986, are gathered in two anthologies edited by Davies and Stone (1995a, 1995b).



Both have useful introductions. The dialog is continued in the papers in Carruthers and Smith 1996. Much of this debate focuses on comparing simulationist approaches to mindreading (as presented in Section 12.5) with the more traditional approach discussed in earlier sections (what is often called the theory theory model of mindreading). Goldman 2006 is a book-length defense of simulationism, written by a philosopher but with extensive discussions of the empirical literature.

Studies on the cognitive neuroscience of mindreading include Apperly et al. 2004, Samson et al. 2004, Samson et al. 2005, Saxe and Kanwisher 2005, Saxe, Carey, and Kanwisher 2004, Tamir and Mitchell 2010, and Waytz and Mitchell 2011. Reviews can be found in Adolphs 2009, Saxe 2009, Carrington and Bailey 2009, Abu-Akel and Shamay-Tsoory 2011, Frith and Frith 2012, and Schurz and Perner 2015. Claims about the modularity of mindreading are critically discussed in Apperly, Samson, and Humphreys 2005. For skepticism about the false photograph task, see Perner and Leekam 2008.

Research into mirror neurons has been reported in many papers – see, for example, Rizzolatti, Fogassi, and Gallese 2001. The findings are presented for a general audience in Rizzolatti, Fogassi, and Gallese 2006 (article) and Rizzolatti and Sinigaglia 2008 (book). For more recent reviews, see Rizzolatti and Sinigaglia 2010 and Rizzolatti and Fogassi 2014. Cook et al. 2014 presents an alternative perspective on mirror neurons, and Kilner and Lemon 2013 offer an independent assessment of the literature.

For more information on empirical findings about emotion recognition in brain-damaged and normal patients, see Adolphs et al. 1994, Phillips et al. 1997, Adolphs and Tranel 2000, and Wicker et al. 2003. See also Croker and Macdonald 2005 and (for a review) Bornhofen and McDonald 2008.





## CHAPTER FIFTEEN

# The Cognitive Science of Consciousness

### OVERVIEW 379

**15.1 The Challenge of Consciousness: The Knowledge Argument** 380

**15.2 Information Processing without Conscious Awareness: Some Basic Data** 382

Consciousness and Priming 382

Nonconscious Processing in Blindsight and Unilateral Spatial Neglect 384

**15.3 So What Is Consciousness For?** 387

What Is Missing in Blindsight and Spatial Neglect 389

Milner and Goodale: Vision for Action and Vision for Perception 389

What Is Missing in Masked Priming 392

**15.4 Two Types of Consciousness and the Hard Problem** 393

**15.5 The Global Workspace Theory of Consciousness** 396

The Building Blocks of Global Workspace Theory 396

The Global Neuronal Workspace Theory 397

**15.6 Conclusion** 400



## Overview

Consciousness is an almost bipolar topic in contemporary cognitive science. On the one hand, we have many exciting experiments and creative theories aiming to understand what consciousness is and how it contributes to cognition. On the other, there are powerful arguments that it is impossible to give an information-processing model of consciousness. This chapter looks at both sides of the debate.

Section 15.1 introduces the challenge of consciousness through Frank Jackson's much-discussed Knowledge Argument. We then consider the differences between conscious and nonconscious information processing. Section 15.2 explores how these are revealed in priming experiments and by studying the behavior of brain-damaged patients. Section 15.3 draws on these findings to explore theories about the function of consciousness. In Section 15.4 we look at two powerful arguments objecting to that whole way of proceeding. According to these arguments,

functional approaches to consciousness cannot help us understand what is truly mysterious about consciousness – at best they can shed light on what are sometimes called the “easy” problems of consciousness. Section 15.5 presents the other side of the coin by reviewing one of the best-established approaches to the functional role of consciousness – the *global workspace theory*.

 15.1

## The Challenge of Consciousness: The Knowledge Argument

We can think about the challenge here through two different perspectives on cognitive agents. The dominant approach within cognitive science has been to look at cognitive agents from the third-person perspective. Cognitive scientists typically work backward from observable behaviors and capacities to information-processing mechanisms that could generate those behaviors and support those capacities. As we have seen in earlier chapters, they do this using a range of experimental techniques and tools, including psychological experiments, functional neuroimaging, and computational modeling. In adopting this third-person perspective, what cognitive scientists do is broadly continuous with what physicists, chemists, and biologists do.

From this third-person perspective what cognitive scientists are working with and trying to explain are publicly observable phenomena – reaction times, levels of blood oxygen, verbal reports, and so forth. But there is another perspective that we have not yet discussed. This is the first-person perspective. Human cognitive agents have sensations. They experience the distinctive smell of a rose, the distinctive sound of chalk on a blackboard, the distinctive feel of cotton against the skin. They react emotionally to events and to each other. They regret the past and have hopes and fears for the future. From the first-person perspective we have a rich, conscious life, full of feelings, emotions, sensations, and experiences. These are all vital parts of what make us human. How can we make sense of them within the information-processing model of the mind?

We can bring some of the problems into clearer focus through a thought experiment originally proposed by the philosopher Frank Jackson. It is usually called the Knowledge Argument. Here is the Knowledge Argument in Jackson’s own words:

Mary is confined to a black-and-white room, is educated through black-and-white books and through lectures relayed on black-and-white television. In this way she knows everything there is to know about the physical nature of the world. She knows all the physical facts about us and our environment, in a wide sense of “physical” which includes everything in *completed* physics, chemistry, and neurophysiology . . .

It seems, however, that Mary does not know all that there is to know. For when she is let out of the black-and-white room or given a color television, she will learn what it is to see something red . . .



After Mary sees her first ripe tomato, she will realize how impoverished her conception of the mental life of *others* has been *all along*. She will realize that there was, all the time she was carrying out her laborious investigations into the neurophysiologies of others, something about these people she was quite unaware of. All along their experiences (or many of them, those got from tomatoes, the sky, . . .) had a feature conspicuous to them, but until now hidden from her.

(Jackson 1986, original emphasis)

When Jackson originally formulated the Knowledge Argument he offered it as a refutation of the philosophical theory known as physicalism (or materialism). According to physicalism, all facts are physical facts. Physicalism must be false, Jackson argued, because in her black-and-white room Mary knew all the physical facts that there are to know and yet there is a fact that she discovers when she leaves the room – the fact about what it is like for someone to see red.



### **Exercise 15.1** State physicalism in your own words. Do you think that Jackson's Knowledge Argument gives a compelling reason to reject physicalism?

Jackson no longer believes that the Knowledge Argument refutes physicalism, however, and so we will not pursue that issue here. For our purposes what is important is that the Knowledge Argument can also be used to argue that information-processing models of the mind are inadequate. The argument would go like this.

- 1 In her black-and-white room Mary has complete knowledge of how information is processed in the brain.
- 2 So, in her black-and-white room Mary knows everything that there is to know about the information processing going on when a person has the experience of seeing red.
- 3 When she leaves the black-and-white room, Mary acquires new knowledge about what goes on when a person has the conscious experience of seeing red.
- 4 Therefore, there must be some aspects of what goes on when a person has the conscious experience of seeing red that cannot be understood in terms of how information is processed in the brain.

The Knowledge Argument raises a powerful challenge to the basic framework assumption of cognitive science that we can give a complete information-processing account of the mind. Little is more salient to each of us than our first-person conscious experience of the world. If, as the Knowledge Argument claims, this is something that cannot be captured in an information-processing account of the mind, then we will have to do a very fundamental rethink of the limits and scope of cognitive science.

For some cognitive scientists, the problem of consciousness is the “last, great frontier.” For many outside the field, in contrast, consciousness reveals the fatal flaw in cognitive science. We will certainly not settle the issue in this book. But the remainder of this chapter surveys some important and exciting research in contemporary cognitive science in the context of this challenge to the very possibility of a cognitive science of consciousness.

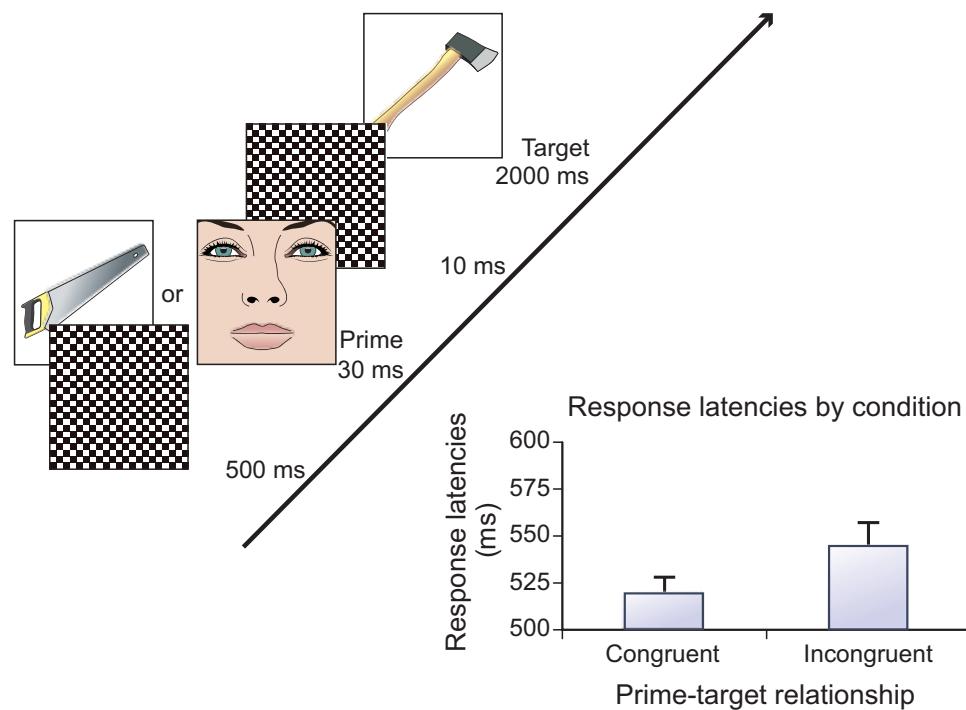

**15.2**

## Information Processing without Conscious Awareness: Some Basic Data

One way to explore the nature of consciousness is to look at the types of information processing and problem solving that can take place without consciousness and compare them with those that seem to require consciousness. Among other things, that will help us understand the function of consciousness – what it actually contributes to cognitive and affective life. We look at two techniques for studying information-processing without consciousness – priming experiments and double dissociations in cognitive neuropsychology.

### Consciousness and Priming

In a typical priming experiment, as illustrated in Figure 15.1, subjects are exposed very briefly to some stimulus – an image on a screen, perhaps, or a sound. The time of exposure is short enough that the subjects do not consciously register the stimulus. Nonetheless, the exposure to the stimulus affects their performance on subsequent tasks – how they



**Figure 15.1** An illustration of a typical priming experiment. The images above the arrow depict the sequence and timing of each stimulus when a tool is the target. The graph shows that people who were presented with a congruent prime were faster to identify the target than people who were presented with an incongruent prime. (From Finkbeiner and Forster 2008)



complete word fragments, for example, or how quickly they can perform a basic classification – processing that can be carried out nonconsciously.

In the experiment in Figure 15.1, subjects are asked to categorize the target as either a face or a tool. There are two different types of prime. One type is congruent with the target (e.g., another tool, if the target is a tool). The other is not congruent (e.g., a tool, if the target is a face). The experiment measures the response latency (the time it takes the subject to classify the target correctly). As the graph illustrates, the experimenters found a significant priming effect for congruent prime-target pairs.

What does this reveal? Think about what counts as a congruent prime-target pair. Figure 15.1 gives one example – a saw and a hammer. These are congruent because they both fall under a single category. Noncongruent prime-target pairs fall under different categories. So, what the priming effect appears to show is that the information processing required to carry out basic categorization can take place nonconsciously. The processing time for correctly classifying a congruent target is less than for a noncongruent target, the standard explanation runs, because the subject is already thinking nonconsciously about the relevant category.

A number of cognitive scientists have objected to the priming paradigm, worried about how one can show that primes really are invisible. A typical method of doing this is to identify a threshold by progressively lowering the presentation time of a stimulus until subjects identify it at chance. This is supposed to show that any stimulus presented for a length of time at or below the threshold will be nonvisible and nonconscious. But one problem with this is that the threshold of visibility can vary. There is some evidence that primes become more visible when they are followed by congruent targets. Varying the mask can also alter the threshold of visibility.

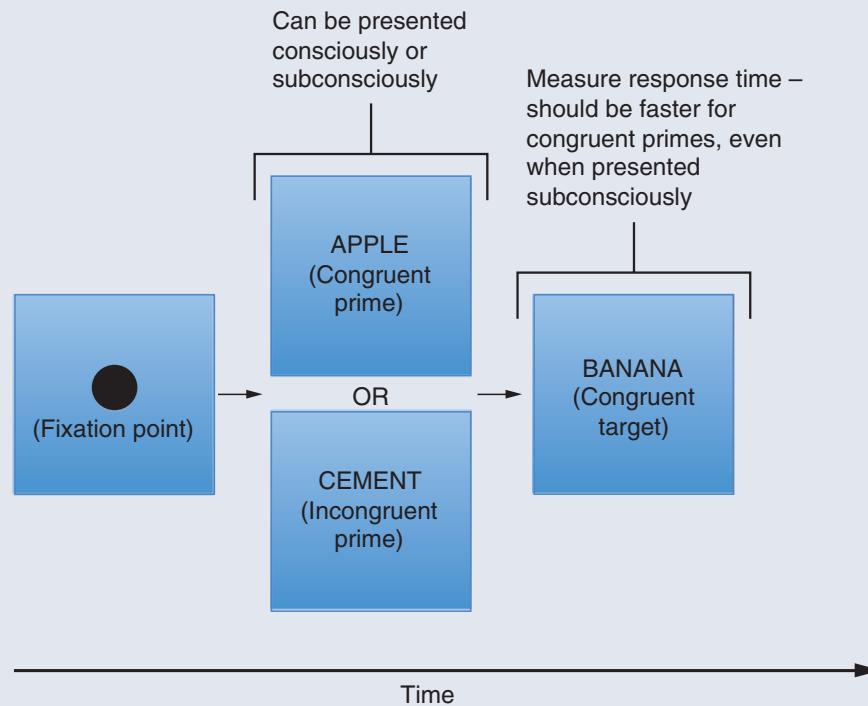
More recent studies have been very sensitive to these methodological issues, and the majority view now is that priming effects do occur, and hence that there are some kinds of nonconscious information processing. The crucial question is: How “smart” is this nonconscious information processing?

One of the most important areas where priming effects have been studied is language processing. In these experiments primes and targets are both words. The most controversial experiments have focused on what is known as *semantic priming* (as illustrated in Box 15.1). Semantic priming occurs when there is a priming effect that can only be explained through information processing about the meaning of words – as opposed, for example, to their phonology (how they are pronounced) or their orthography (how they are spelled).

There is interesting evidence for semantic priming from studies with bilingual subjects where prime and target are in different languages, particularly where those languages are in very different scripts (Chinese and English, for example). Many studies have shown robust priming effects when subjects are asked to decide whether or not a target string of letters is a proper word or not (what is called the *lexical decision task*). Interestingly, the priming effect tends to occur only when the prime is in the dominant (first) language (L1) and the target is in the second language (L2).

Semantic priming is potentially very significant, because semantic processing is widely held to be very high-level and dependent upon conscious awareness. To go back to the

### BOX 15.1 A Typical Semantic Priming Experiment



There are many variations. Sometimes the words are presented in different languages, as discussed in the main text, and sometimes the semantic congruence varies for the target instead of the prime. Participants can be asked to hit a button simply when they see the target word or make some more difficult judgment about the word (e.g., whether it is in fact a word).

distinction we looked at in Chapter 8, semantic processing has typically been thought to be nonmodular (as opposed to processes such as phonological parsing, often thought to be modular). So, semantic priming is important because it seems to show that there can be information processing that is both nonmodular and nonconscious.



**Exercise 15.2** Explain in your own words why the distinction between modular and nonmodular information processing is important for thinking about nonconscious information processing.

### Nonconscious Processing in Blindsight and Unilateral Spatial Neglect

Cognitive neuropsychologists study cognitive disorders, primarily resulting from brain damage. The guiding idea is that we can work backward from what happens when things



go wrong to how they function in the normal case. So, for example, if in one type of brain damage we see ability A functioning more or less normally while ability B is severely impaired, then we can infer that in some sense A and B are independent of each other – or, as cognitive neuropsychologists call it, we can infer a *dissociation* between them.

A *double dissociation* occurs when we have a dissociation in each direction – that is, in one disorder we have ability A functioning normally with B significantly impaired, while in a second disorder we have ability B functioning normally with A significantly impaired. Double dissociations provide stronger evidence that A and B are independent of each other.



### Exercise 15.3 Explain in your own words why a double dissociation is a better sign of independence than a single dissociation.

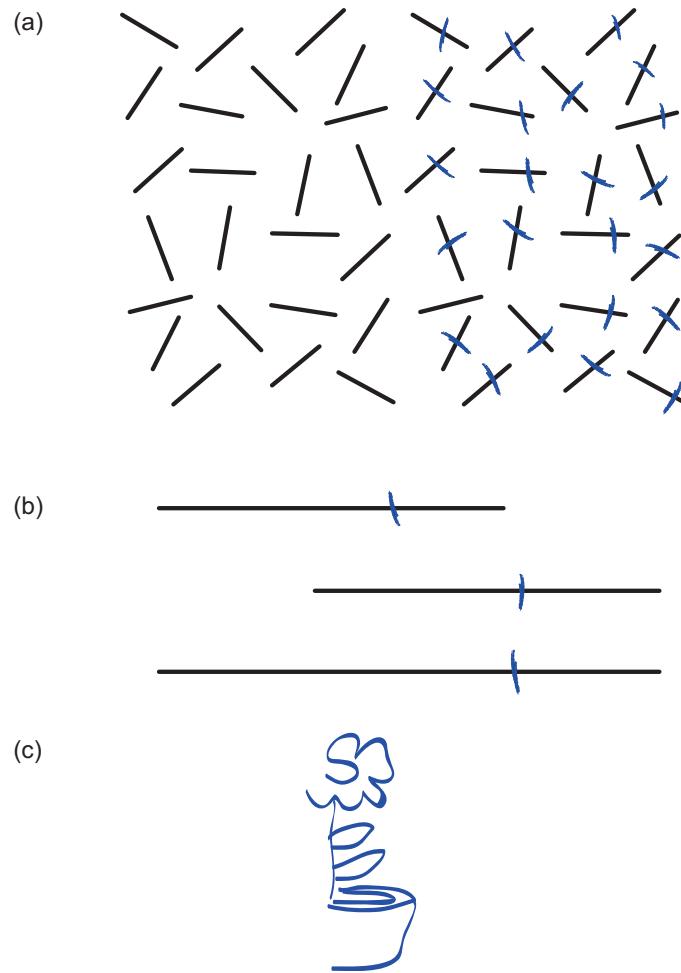
Cognitive psychologists studying psychological disorders caused by brain trauma have identified very interesting dissociations involving consciousness. There are surprisingly many tasks that can be carried out nonconsciously by brain-damaged patients, even though they are typically performed with conscious awareness by normal subjects. In particular we will look at two much-studied disorders – unilateral spatial neglect and blindsight.

Unilateral spatial neglect (also known as *hemiagnosia* or *hemineglect*) is relatively common and typically occurs after damage to the right hemisphere, particularly damage to the parietal and frontal lobes. Its defining feature is that patients lack awareness of sensory events on the contralateral side of space (on the opposite side of the world to the side of the brain that is damaged). In the vast majority of cases, the neglected side is the left-hand side.

The neglect phenomenon was very strikingly illustrated by two Italian neuropsychologists in 1978. Eduardo Bisiach and Claudio Luzzatti asked two neglect patients to describe from memory the central square in Milan with the famous Duomo (cathedral). The patients were initially asked to describe the square as if they were standing in front of the Duomo. As predicted, the patients failed to describe the houses and shops on the left-hand side of the square (from their vantage point in front of the Duomo). Bisiach and Luzzatti then asked the patients to orient themselves differently, so that they were imagining themselves on the edge of the square looking at the Duomo. Now the patients accurately described the houses and shops they had previously neglected, and instead missed out the side of the square that they had previously described. Figure 15.2 shows further examples of typical visual deficits in neglect patients.

Neglect also affects action. A neglect patient might only shave or apply makeup to one side of their face, for example. Or they might eat only from one side of a plate.

The blindsight patients who have been most studied report little to no awareness in one side of their visual field. They have what is called a *scotoma* (a region of very diminished visual acuity that does not occupy the whole visual field). In both blindsight and unilateral spatial neglect, patients report themselves to be unaware of what is going on in part of their visual field. The etiology (cause) is different, however. The impairment in blindsight is typically due to lesions in the primary visual cortex (V1, or the striate cortex).



**Figure 15.2** Examples of deficits found in patients with left spatial neglect (damage to the right hemisphere of the brain). (a) Unilateral neglect patients typically fail to mark the lines on the contralesional (here, left) side of a sheet of paper. (b) Patients are asked to bisect each line. Their markings are typically skewed to the right, as if they do not see the leftmost segment. (c) Patients are either asked to draw something from memory or to copy another illustration placed in front of them. In both cases, unilateral neglect patients tend to omit parts on the contralesional side. (From Driver and Vuilleumier 2001)

For our purposes, the interesting feature of both blindsight and unilateral spatial neglect is that patients appear to have surprising residual visual functioning despite reporting a more or less complete lack of visual awareness. Blindsight patients can respond to stimuli in the scotoma, and visual neglect patients can respond to stimuli in the neglected region of space.

One challenge in exploring the residual abilities of blindsight patients is that they will often find the experiments absurd. Ernst Pöppel, whose important 1973 article coauthored with Douglas Frost and Richard Held was one of the first to study blindsight, reported a patient irritatedly saying “How can I look at something that I haven’t seen?” when asked to direct his eyes to a target in his blind field.



To overcome this challenge, experimenters have used nonverbal forced choice tests. In essence, patients are forced to guess in situations where they feel that they have no basis to make a judgment or to perform an action. The choices are usually binary – is the stimulus moving or stationary, is it high or low in the visual field, is it horizontal or vertical?

Experimenters often find that blindsight patients perform significantly better than chance, even when the patients describe themselves as guessing (and so would be expected to perform at chance levels). There is strong evidence that blindsight patients can localize unseen stimuli in the blind field, that they can discriminate orientation, and that they can detect moving and stationary figures randomly interspersed with blank trials.

Neuropsychologists have also found that blindsight patients are capable of some types of form perception. Here is an example from a striking set of experiments performed by Ceri Trevethan, Arah Sahraie, and blindsight pioneer Larry Weiskrantz, working with a patient known by his initials D.B. Figure 15.3 depicts line drawings of animals that were presented within D.B.'s blind field.

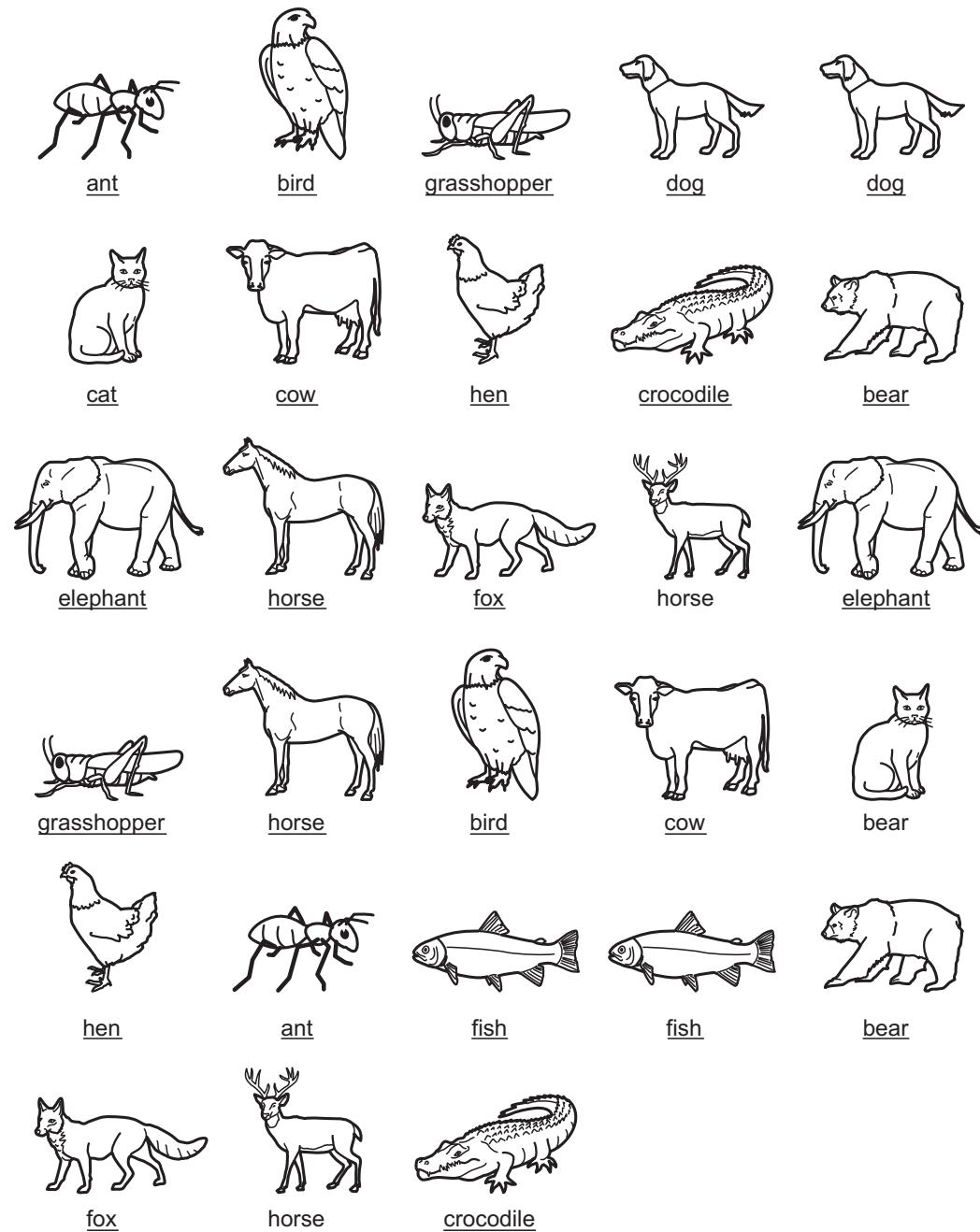
The figures were shown at very low contrast (2 percent – although they are depicted in high contrast in Figure 15.3). The patient was told that he was being shown a picture of an animal and asked to guess which animal it was. The figure indicates the responses given, with correct answers underlined. As illustrated, D.B. achieved 89 percent accuracy, despite reporting no awareness whatsoever of any of the figures.

Spatial neglect patients also have considerable residual abilities. A famous example identified by neuropsychologists John Marshall and Peter Halligan is illustrated in Figure 15.4. Marshall and Halligan showed P.S., a neglect patient, the two pictures in the diagram – one of a normal house and one of a house on fire. Since the flames were on the left-hand side of the picture, P.S. did not report seeing any difference between the two pictures. Nonetheless, when asked which house she would prefer to live in, P.S. reliably chose the house that was not on fire (9 times out of 11).

Italian neuropsychologists Anna Berti and Giacomo Rizzolatti used a semantic priming paradigm to explore whether neglect patients could identify semantic categories in their neglected field. Neglect patients were presented with priming stimuli in their neglected visual field and then asked to categorize objects presented in the normal field. As discussed above, the assumption for priming experiments is that, when the prime stimulus and the target stimulus are congruent (i.e., from the same category), then categorization will be easier and quicker, provided that the prime stimulus is processed. Berti and Rizzolatti found the predicted effect in patients who denied all awareness of the prime stimuli and so concluded that semantic information is processed in the neglected visual field.

## 15.3 So What Is Consciousness For?

Experiments on both brain-damaged and normal subjects indicate that many information-processing tasks can be performed without conscious awareness – or, more precisely, can be



**Figure 15.3** D.B.'s responses to pictures of animals presented in his blind field. Correct answers are underlined. (From Trevethan, Sahraie, and Weiskrantz 2007)

performed by subjects who do not report any conscious awareness of the discriminations and selections that they are making. This leaves us with a puzzle. What exactly does consciousness contribute? Why do we need it? In order to make progress on this we need to look, not just at what blindsight and neglect patients can do, but also at what they can't do.



**Figure 15.4** An illustration of the two houses presented to P.S. The houses are identical, except that one has flames shooting out of its left side. Because P.S. possesses left-side spatial neglect, she reported not being able to see the flames but still consistently selected the other house when asked which house she would prefer to live in. (From Marshall and Halligan 1988)

## What Is Missing in Blindsight and Spatial Neglect

Blindsight and neglect patients have considerable residual abilities that can be teased out with careful experiments, but there are still massive differences between these patients and normal subjects.

One very striking fact is just how difficult it is to elicit the residual abilities. Neither blindsight nor neglect patients will voluntarily do things in their blind or neglected fields. This is most obvious in neglect patients. What characterizes the disorder is not just that patients report a complete lack of awareness of what is going on in the neglected visual field. It is also that they do not direct any actions within those regions of space that fall within the neglected visual field. This is the case both for their own personal, bodily space (so that male patients do not shave on the neglected side of their face) and for external space (so that they do not direct actions at objects located on the neglected side of the world as they perceive it). The same holds for blindsight patients, who never initiate actions toward the blind field, despite being able to point to stimuli in the blind field (when forced to do so).

This suggests a hypothesis about the difference between conscious and nonconscious information processing. Both normal and brain-damaged subjects receive many different types of nonconscious information about the world and about their own bodies. But subjects can only initiate voluntary actions on the basis of information that is conscious. Only conscious information allows subjects to identify targets and to plan actions toward them. The neuropsychologists David Milner and Melvyn Goodale have developed a sophisticated theory of vision that is built around this idea that one of the roles of consciousness is to permit voluntary and deliberate action.

## Milner and Goodale: Vision for Action and Vision for Perception

Milner and Goodale's theory is based on the existence of two anatomical pathways carrying visual information in the primate brain. We looked at some of the neurophysiological

evidence for these two anatomical pathways in Section 3.2 when we reviewed the important Mishkin and Ungerleider experiments. Visual information takes two different routes from the primary visual cortex. One pathway, the ventral pathway, projects to the temporal lobe. A second pathway, the dorsal pathway, carries information to the posterior parietal lobe. (See Figure 3.5 for an illustration of the two pathways.)

The two pathways have very different functions. For Mishkin and Ungerleider, as we saw in Chapter 3, the crucial functional distinction is between the “what” system, concerned with object identification and subserved by the ventral pathway, and the “where” system, concerned with locating objects in space. Milner and Goodale have a related but somewhat different interpretation. They distinguish two types of vision, which they term *vision for action* and *vision for perception*.

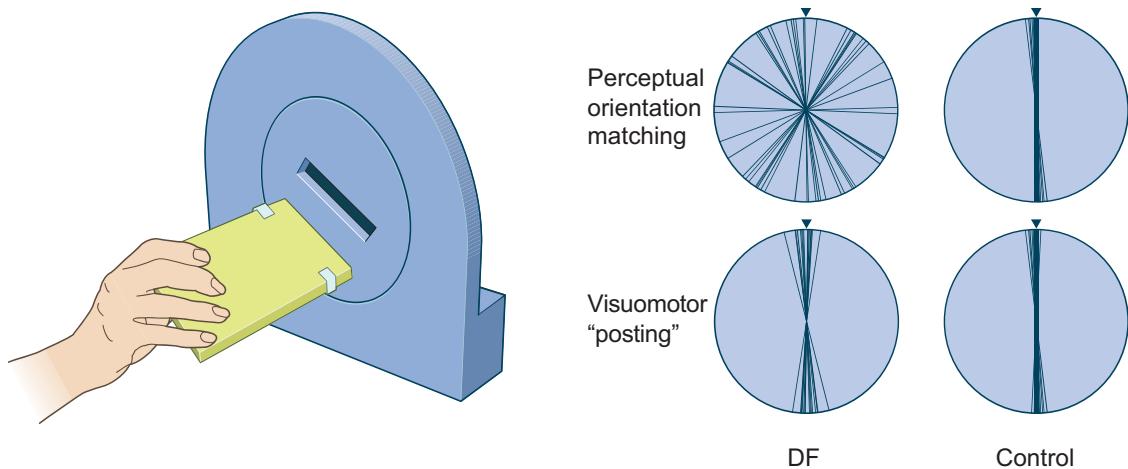
Vision for action has to do with how actions are executed. When you reach out to grasp something, for example, your hand automatically prepares itself, so that the fingers are at an appropriate aperture. This involves complex information processing, including estimates of the likely size of the object, taking into account distance and so forth. We all also constantly engage in online correction of movement, compensating for environmental change or initial errors of trajectory. The relevant processing here takes place in the dorsal pathway;

Vision for perception, on the other hand, deals with actually initiating deliberate action. Vision for perception allows targets to be identified and goals to be set. It depends upon information from the ventral stream. According to Goodale and Milner, only information relevant to what they call vision for action is actually conscious. Conscious awareness is restricted to the ventral pathway while the dorsal stream governs the visual control of movement nonconsciously.

Milner and Goodale rely heavily on experimental studies of both normal and brain-damaged patients. Here are two examples that illustrate how consciousness is and is not involved in vision.

Milner and Goodale’s patient D.F. is one of the most studied and important neuropsychological patients. After carbon monoxide inhalation, D.F. developed what is known as *visual form agnosia*, substantially impaired visual perception of shape and orientation. The neural damage underlying her agnosia involved very serious damage to the ventral pathway (i.e., the system responsible for consciously selecting goals and initiating action).

D.F. performs many visuomotor tasks successfully, even though she is unable to recognize or identify the relevant features in her environment. Figure 15.5 illustrates a much-discussed example of two tasks where D.F. performs very differently. When asked to “post” a card into a slot D.F. was able to match her movements to the orientation of the slot and performed almost as successfully as normal subjects. But when asked to make an explicit judgment about the slot’s orientation D.F.’s responses were almost random. This was the case whether she was asked to describe the orientation verbally or nonverbally (by rotating a card to match the orientation). According to Milner and Goodale, D.F. is receiving nonconscious information about orientation through the dorsal pathway, but because of damage to her ventral pathway is not consciously aware of the orientation.

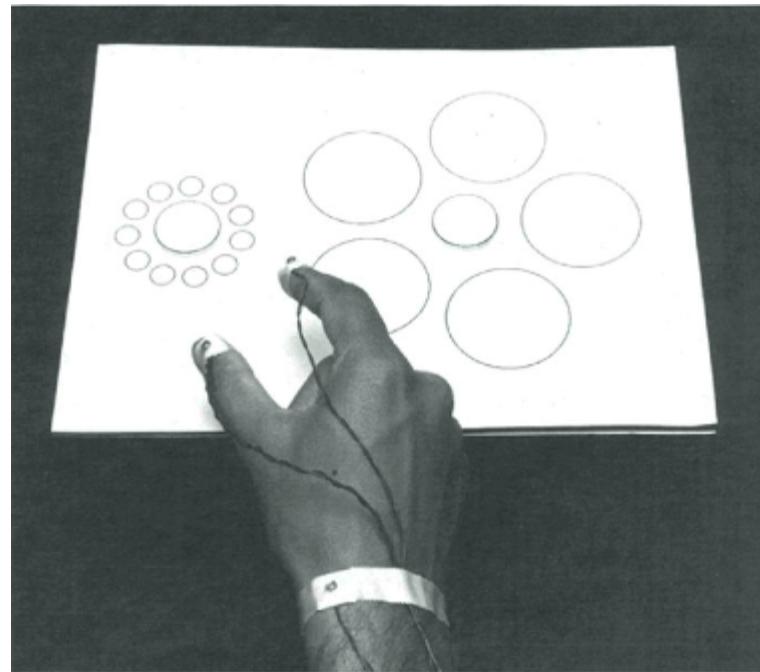


**Figure 15.5** In this experiment, subjects were asked either to “post” a card into a slot or to rotate another hand-held card to match the orientation of the slot. The angle of the slot varied across trials, although in each case, the diagrams have been normalized so that the correct result is vertical. Normal subjects can perform both tasks with little difficulty. Patient D.F., in contrast, can carry out the visuomotor task almost as well as normal subjects, but her responses in the explicit matching task are almost random. (From Milner and Goodale 1998)

Visual illusions provide another source of evidence for the dissociation between (nonconscious) vision for action and (conscious) vision for perception. Visual illusions affect how subjects consciously perceive the size and shape of objects. A number of experimenters have found, however, that the illusion does not carry over to visuomotor behavior. Subjects will report seeing an illusion, but when asked to make appropriate movements they will configure their grip and make other adjustments according to the correct dimensions of the relevant objects, not the dimensions that they report perceiving. So conscious perception (vision for perception) dissociates from (nonconscious) information relevant to the control of visuomotor behavior (vision for action). Figure 15.6 illustrates the experiment used by Aglioti, DeSouza, and Goodale to identify this dissociation, utilizing what is known as the Ebbinghaus illusion.

In addition, neuroimaging studies, such as those published by Fang Fang and Shen He in 2005, suggest that ventral stream activity is correlated with consciousness, while activity in the dorsal stream is not. Fang and He compared activation levels in areas known to be very involved in object processing in the dorsal and ventral streams, respectively. They used a technique known as *interocular suppression* in which one eye is presented with an image of an object while the other eye is presented simultaneously with a high-contrast pattern that blocks conscious awareness of the presented image.

This paradigm enabled Fang and He to examine activation levels in the dorsal and ventral streams in the absence of conscious awareness and to compare those levels with activation levels when conscious awareness of the image was not suppressed. They found robust levels of activity in the dorsal stream even in the nonconscious conditions. In contrast, ventral stream activation was confined to the conscious condition.



**Figure 15.6** In the Ebbinghaus illusion, two circles are illusorily seen as differently sized, depending on what surrounds them. The figure illustrates experiments published by Aglioti, DeSouza, and Goodale in 1995. The experimenters measured the size of the opening between fingers and thumb when subjects were asked to pick up two disks that they reported as being differently sized. They found no significant differences in grip aperture, suggesting that this aspect of the fine-grained control of grasping draws on different types of visual information than those that yield conscious awareness of the disks.

In conclusion, Milner and Goodale's distinction between (conscious) vision for perception and (nonconscious) vision for action, together with the evidence supporting it from brain-damaged and normal subjects, both supports and clarifies the hypothesis that consciousness is important for initiating action. If Milner and Goodale are correct, then conscious awareness is key for identifying targets and for macro-level planning for how to effect actions. But conscious awareness is not typically involved in the fine-grained, online control of bodily movements.

## What Is Missing in Masked Priming

Masked priming experiments offer powerful evidence of nonconscious semantic processing. At the same time, though, information-processing in priming experiments is very different from normal, conscious information-processing.

One key finding here is it is hard to retain information without conscious awareness. So, semantic information processed below the threshold of consciousness in masked priming experiments is very transitory and short-lived. Here is an illustration from experiments published by Anthony Greenwald, Sean Draine, and Richard Abrams in 1996. The authors



used a typical categorization task, asking subjects to identify first names as male or female or to classify words as pleasant or unpleasant in meaning. They succeeded in eliciting a robust priming effect when subjects were presented with a congruent masked prime. This effect was present both when the prime was presented subliminally and when it was presented supraliminally (above the threshold of consciousness). This allowed Greenwald, Draine, and Abrams to study the differences between subliminal priming and supraliminal priming. The particular dimension they explored was what happened when they varied the time between prime and trial (the so-called *stimulus-onset asynchrony*, SOA).

Greenwald, Draine, and Abrams found a significant difference. In supraliminal cases, where the subjects were conscious of the prime, the priming effect was robust across all SOAs. The length of the delay between prime and target did not make a significant difference. In contrast, in the subliminal cases, with the subjects not consciously perceiving the prime, the effect was robust only at the shortest intervals and disappeared completely once the SOA went above 100 ms.

This experiment suggests an additional hypothesis about the function of conscious awareness, namely, that consciousness allows information to be explicitly retained and maintained. According to this hypothesis, information that is picked up nonconsciously can indeed be deployed in relatively sophisticated tasks, but it can be used only within a very limited time horizon. Conscious information, in contrast, is more transferable and flexible. It can be used beyond the here and now. There are definite parallels between this idea and Goodale and Milner's distinction. Vision for action is restricted to the online control and fine-tuning of behavior. It does not persist in the way that conscious visual information persists. That is one reason why Goodale and Milner think that the conscious vision-for-perception system is required for high-level action-planning.

## 15.4

### Two Types of Consciousness and the Hard Problem

Two related ideas have emerged about the function of consciousness. First, conscious awareness seems extremely important for planning and initiating action (as opposed to the online control of behavior, which can be carried out through nonconscious information processing). Second, conscious information persists longer than nonconscious information. In the next section we will look at one example of a theory of consciousness that can accommodate these two ideas. First, though, we need to consider some important concerns about this whole way of proceeding that have been raised by the philosophers Ned Block and David Chalmers.

The philosopher Ned Block has cautioned cognitive scientists to be very careful about drawing conclusions about the nature and function of consciousness from neuropsychological disorders such as blindsight and unilateral spatial neglect. He thinks that these conclusions rest on flawed inferences. What causes the problem, according to Block, is a confusion between two very different concepts of consciousness. He calls these *phenomenal consciousness* and *access consciousness*. Here is how he characterizes the two notions, which he terms P-consciousness and A-consciousness, respectively:

### *Phenomenal consciousness*

P-consciousness is experience . . . We have P-conscious states when we see, hear, smell, taste, and have pains. P-conscious properties include the experiential properties of sensations, feelings, and perceptions, but I would also include thoughts, wants, and emotions. (Block 1995b)

### *Access consciousness*

A state is A-conscious if it is poised for direct control of thought and action. To add more detail, a representation is A-conscious if it is poised for free use in reasoning and for direct “rational” control of action and speech. (The rational is meant to rule out the kind of control that obtains in blindsight.) (Block 1995b)



### **Exercise 15.4** Give your own examples of A-consciousness and P-consciousness and describe the difference between them in your own words.

From Block’s perspective, the real problem of consciousness is the problem of understanding P-consciousness. All of the things that we have been looking at in the previous section, however, are really examples of A-consciousness. This is the “confusion” that he identifies in the title of his influential paper “On a confusion about a function of consciousness.”

According to Block, the experiments and studies discussed in the previous section ultimately only inform us directly about the function of A-consciousness. They do not directly address the function of P-consciousness. Our two hypotheses about the function of consciousness are hypotheses about the difference between conscious information processing and nonconscious information processing. This does not get to the heart of what Block sees as the real problem of consciousness, which has to do with how and why we experience the world the way we do.

Block’s distinction between A-consciousness and P-consciousness is related to further distinctions drawn by the philosopher David Chalmers in his influential book *The Conscious Mind* and other writings. Chalmers thinks that there is no single problem of consciousness. Instead, he thinks that we need to make a distinction between a cluster of relatively easy problems and a single, really difficult problem – what he calls the “hard problem” of consciousness.

Here are some examples of what Chalmers provocatively identifies as easy problems of consciousness:

- explaining an organism’s ability to discriminate, categorize, and react to environmental stimuli;
- explaining how a cognitive system integrates information;
- explaining how and why mental states are reportable;
- explaining how a cognitive system can access its own internal states;
- explaining how attention gets focused;
- explaining the deliberate control of behavior;
- explaining the difference between wakefulness and sleep.



In Block's terminology these are different aspects of understanding A-consciousness. In the last analysis, they are all problems to do with how an organism accesses and deploys information.

Chalmers recognizes that "easy" is a relative term. None of the so-called easy problems has yet been solved, or even partially solved. The reason he calls them easy problems is that at least we have some idea of what a solution would look like. The easy problems are all problems that are recognizable within the basic framework of cognitive science and scientific psychology. People write papers about them, reporting relevant experiments and constructing theories.

According to Chalmers, though, no amount of progress on the easy problems of consciousness will help with the hard problem. Here is how he characterizes the hard problem:

The really hard problem of consciousness is the problem of *experience*. When we think and perceive, there is a whir of information-processing, but there is also a subjective aspect. As Nagel (1974) has put it, there is *something it is like* to be a conscious organism. This subjective aspect is experience. When we see, for example, we *experience* visual sensations: the felt quality of redness, the experience of dark and light, the quality of depth in a visual field. Other experiences go along with perception in different modalities: the sound of a clarinet, the smell of mothballs. Then there are bodily sensations, from pains to orgasms; mental images that are conjured up internally; the felt quality of emotion, and the experience of a stream of conscious thought. What unites all of these states is that there is something it is like to be in them. All of them are states of experience.

It is undeniable that some organisms are subjects of experience. But the question of how it is that these systems are subjects of experience is perplexing. Why is it that when our cognitive systems engage in visual and auditory information-processing, we have visual or auditory experience: the quality of deep blue, the sensation of middle C? How can we explain why there is something it is like to entertain a mental image, or to experience an emotion? ...

If any problem qualifies as *the* problem of consciousness, it is this one.



### Exercise 15.5 In your own words characterize the proposed distinction between what Chalmers calls the easy problems of consciousness and what he calls the hard problem.

In Chalmers's phrase, looking at what happens in masked priming experiments or at the differences between normal subjects and blindsight patients can only help with the easy problems of consciousness. None of these things can possibly help with the hard problem of consciousness. The differences between normal subjects and patients suffering from blindsight or spatial neglect, or between subliminal and supraliminal, are differences in access to information. They cannot help us understand the nature of experience or what it is to be phenomenally conscious. In fact, Chalmers draws a very drastic conclusion from his distinction between easy and hard problems. He thinks that the hard problem of consciousness is in principle intractable to cognitive science (or any other kind of science).

Without trying to settle the matter one way or the other, it seems plausible that progress is going to depend upon having a better idea of what an information-processing account of access consciousness might look like. Discussing the limits that there might or might not be to a particular type of explanation will be much easier when there is a particular example on which to focus. In the next section we will look at the *global workspace theory of consciousness*, which is an interesting candidate for an information-processing solution to some of the problems that Chalmers identifies as the easy problems of consciousness.

## 15.5

## The Global Workspace Theory of Consciousness

Global workspace theory was originally proposed by the psychologist and cognitive scientist Bernard Baars in his book *A Cognitive Theory of Consciousness*, published in 1988. Since then it has been taken up and developed by many others, including the neuroscientists Antonio Damasio and Stanislas Dehaene, as well as the philosopher Peter Carruthers. More recent presentations (in line with the general turn toward the brain in cognitive science) have emphasized the neural dimension of global workspace theory.

Global workspace is not, of course, the only theory of consciousness currently being discussed by cognitive scientists. But it fits very naturally with many of the topics that we have been discussing in this chapter (and indeed throughout the book). In Block's terminology, global workspace theory is a theory of access consciousness – a theory of how information is made available for high-level cognition, action-planning, and speech. The theory is based on an analysis of the function of consciousness that directly addresses many of what Chalmers identifies as “easy” problems of consciousness. And finally, it draws on ideas that we have discussed earlier in the book – including the idea that the mind has both modular and nonmodular components and the idea that attention serves a “gatekeeper” function in controlling what crosses the threshold of conscious awareness.

## The Building Blocks of Global Workspace Theory

We will focus on the version of global workspace theory presented by Stanislas Dehaene and collaborators. Stanislas Dehaene and Lionel Naccache give a very clear account of the theoretical underpinnings of the global workspace theory in their 2001 article “Towards a cognitive neuroscience of consciousness: Basic evidence and a workspace framework.” They propose the theory as the best way of making sense of the basic functional benefits of consciousness within a framework set by some widely accepted assumptions about the architecture of the mind.

Dehaene and Naccache focus in particular on three different things that they believe consciousness makes possible. These are:

- the intentional control of action
- durable and explicit information maintenance
- the ability to plan new tasks through combining mental operations in novel ways



They consider these three functions of consciousness relative to two basic theoretical postulates about mental architecture and the large-scale organization of the mind.

The first is a version of the modularity theory that we explored at length in Chapter 8. Modular processes have two key features. They are *domain-specific* and *informationally encapsulated*. That is to say, they are each dedicated to solving circumscribed types of problem that arise in very specific areas and in solving those problems they typically work with restricted databases of specialized information. Many cognitive tasks involve a series of modules – executing an action is a good example – but, according to the classical version of modularity theory, there are some cognitive tasks that cannot be carried out by modular systems. These are tasks that are domain-general (they span a range of cognitive domains) and that can only be solved by drawing upon the full range of information that the organism has available to it. The global workspace is in essence a metaphorical name for this type of domain-general information processing.



### Exercise 15.6 Review the discussion of modularity in Chapter 8.

Dehaene and Naccache suggest that the distinction between the conscious and non-conscious minds maps onto the distinction between modular processing and nonmodular processing. Consciousness is restricted to information within the global workspace.

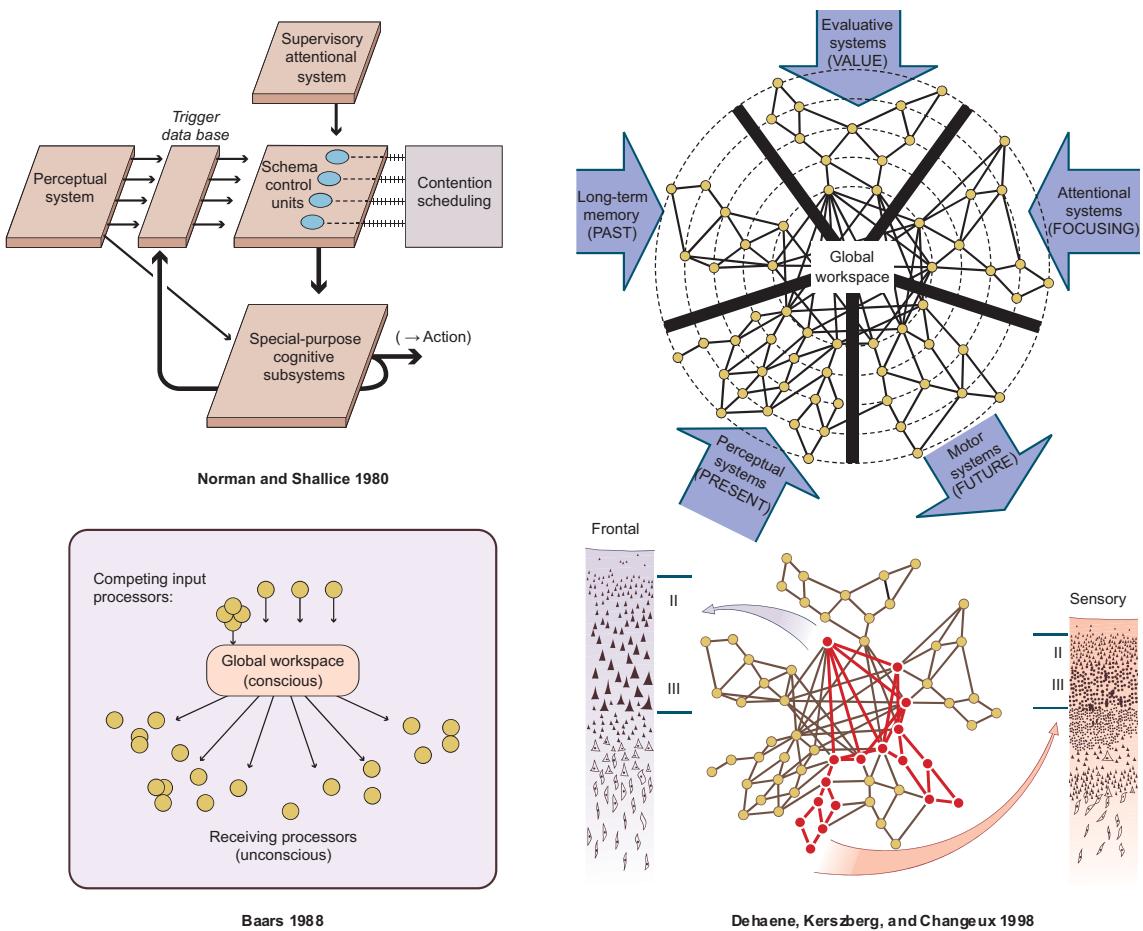
Their second theoretical postulate has to do with how information becomes available to the global workspace. Attention is the key mechanism here. It functions as a gatekeeper, allowing the results of modular information processing to enter the global workspace. For the global workspace theory, attention and consciousness are very closely linked. This way of thinking about the role of attention has a long pedigree within cognitive science, going back to the pioneering work of Donald Broadbent, reviewed in Section 1.4. Attention is thought of both as a *filter* (screening out unnecessary information, as in the cocktail party effect) and as an *amplifier* (allowing information that would otherwise have been unconscious to become available to consciousness).

## The Global Neuronal Workspace Theory

Three versions of the global workspace theory are illustrated in Figure 15.7, showing how the workspace idea has evolved over the last 30 years.

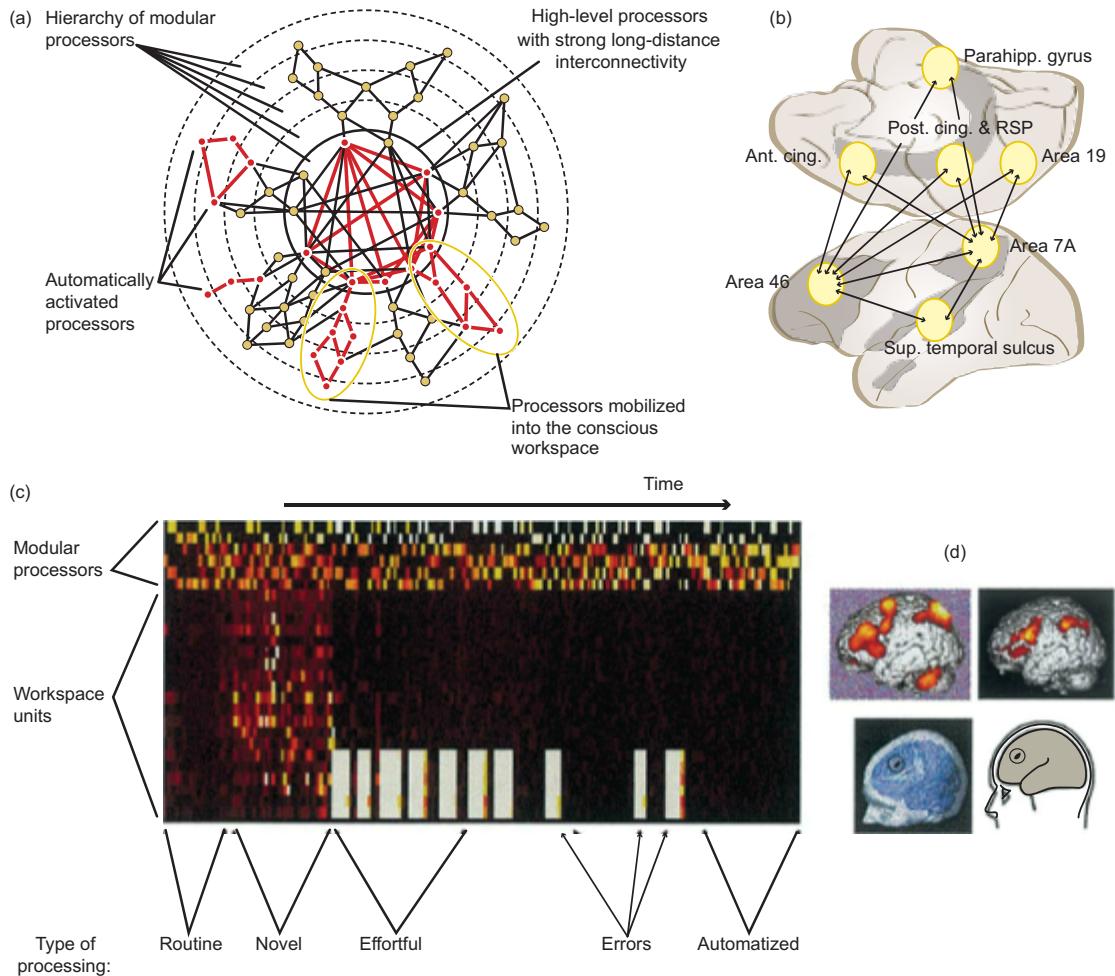
The first version (under a different name) was the theory of attention originally proposed by Donald Norman and Tim Shallice. As the figure illustrates, attention performs what Norman and Shallice term *contention scheduling*. Contention scheduling is required when different cognitive systems propose competing responses (whether cognitive or behavioral) to a single set of stimuli. Contention scheduling effectively resolves the competition to select a single response, which can either be an output to the action systems or can be fed back into the cognitive systems.

The terminology of global workspace was introduced by Bernard Baars in the late 1980s. One version of his theory is depicted in the figure, showing very clearly how the global workspace is envisioned as a conscious window between nonconscious inputs and conscious outputs.



**Figure 15.7** In the Norman and Shallice 1980 model (top left), conscious processing is involved in the supervisory attentional regulation, by prefrontal cortices, of lower-level sensorimotor chains. According to Baars 1988, conscious access occurs once information gains access to a global workspace (bottom left), which broadcasts it to many other processors. The global neuronal workspace (GNW) hypothesis (right) proposes that associative perceptual, motor, attention, memory, and value areas interconnect to form a higher-level unified space where information is broadly shared and broadcasted back to lower-level processors. The GNW is characterized by its massive connectivity, made possible by thick layers II/III with large pyramidal cells sending long-distance cortico-cortical axons, particularly dense in prefrontal cortex. (From Dehaene and Changeux 2011)

A much more recent version of the theory is depicted on the right side of Figure 15.7. It was developed by Stanislas Dehaene, Michel Kerszberg, and Jean-Pierre Changeux. This shares some features with the other two versions, particularly the idea that the global workspace receives inputs from different cognitive modules and then sends outputs to motor systems. The Dehaene, Kerszberg, and Changeux theory is much more strongly grounded in hypotheses about neural implementation and connectivity – which is why they call their theoretical construct the *global neuronal workspace*, rather than simply the global workspace. This emerges even more clearly in Figure 15.8.



**Figure 15.8** The neural substrates of the global workspace. (a) Hierarchy of connections between different processors in the brain. Note the strong long-distance connections possessed by the higher levels. (b) Proposed anatomical substrate of the global workspace. This includes a network linking the dorsolateral prefrontal, parietal, temporal, and anterior cingulate areas with other subcortical regions (RSP = retrosplenial region). (c) Neural dynamics of the global workspace, derived from a neural simulation of the model shown in (a). The activation levels of various processor units (top lines) and workspace units (bottom lines) are shown as a function of time. (d) Different parts of the global workspace network activated by different tasks, including generation of a novel sequence of random numbers, effortful arithmetic, and error processing. (From Dehaene and Naccache 2001)

Figure 15.8 makes plain the distributed nature of the global neuronal workspace, as envisaged by Dehaene and his collaborators. They see the modular part of the mind as composed of many interconnecting modules that feed into each other in a hierarchical manner. (The hierarchy is depicted by concentric circles, and the closer the circles to the center the higher their place in the hierarchy.) Some of the hierarchical modules form fully automatic and nonconscious networks. Others in contrast have amplified levels of activity that allow them to feed into the global workspace.

The global neuronal workspace itself is not a single neural location (as the metaphor of a workspace might initially suggest), but rather a distributed network of high-level processors that are highly connected to other high-level processes. The candidate areas identified include the prefrontal, parieto-temporal, and cingulate cortices – all areas that we have discussed in the context of different types of high-level cognition at various points in this book. The lower portion of Figure 15.8 includes a neural network simulation of the global neuronal workspace and when it becomes engaged, in addition to an fMRI diagram indicating activation levels across the hypothesized network during high-level conscious tasks such as mental arithmetic.

The hypothesis is that the global neuronal workspace is generated by the activities of a particular type of neuron called pyramidal neurons. Pyramidal neurons are very widespread in the mammalian cortex and particularly dense in the prefrontal, cingulate, and parietal regions (all hypothesized to be important in the global neuronal workspace). They are characterized by a single long axon and heavily branched dendrites, which allow them to communicate with many other neurons and with distant brain areas. Dehaene and collaborators hypothesize that networks of pyramidal neurons connect specialized modular processes and allow their outputs to be broadcast across the brain so that they are available for action-planning, verbal report, and other high-level cognitive processes.



**Exercise 15.7** Explain in your own words how the global neuronal workspace theory incorporates the hypotheses about the function of consciousness identified in Section 14.4.

## 15.6 Conclusion

There are two very different approaches to consciousness. On the one hand there are those who think that consciousness is a mystery that we have no idea how to tackle with the tools and methods of cognitive science. On the other hand, we have thriving research programs that study different aspects of the conscious mind and how consciousness contributes to action and cognition.

The “mysterians,” as they are sometimes called, hold that the various research programs we have looked at only touch upon the “easy” aspects of the problem of consciousness – at best they can only explain access consciousness, as opposed to the really tough problem of explaining how and why we are phenomenally conscious. The global neuronal workspace theory was primarily developed to explain how consciousness can make a difference to cognition. The theory gives an account of why some information becomes conscious and how that information has a distinctive role to play in higher-level cognition. Mysterians will say that this account is all well and good, but cannot come to grips with the “hard problem” of explaining the distinctive experience of being conscious.



In response, cognitive scientists working on consciousness may well respond that the so-called hard problem of consciousness will disappear once we have a good enough understanding of the various phenomena lumped together under the label “access consciousness.” This is the view taken by the philosopher Daniel Dennett, whose books *Content and Consciousness* and *Consciousness Explained* have been very influential in discussions of consciousness.

Perhaps the source of the problem is that we do not have any real idea of what a complete account of access consciousness would look like. As its originators would be the first to admit, the global neuronal workspace theory is programmatic in the extreme. So are its competitors. It may well be that if we were in possession of something much more like a complete theory then we would be less likely to have intuitions about the impossibility of solving the so-called hard problem. What makes the intuitions seem compelling is that our knowledge is so incomplete and our investigations of the cognitive science of consciousness at such an early stage.

There is a lesson to be learned from debates in the nineteenth and early twentieth century about vitalism in biology. Vitalists such as the philosopher Henri Bergson and the biologist John Scott Haldane believed that the mechanist tools of biology and chemistry were in principle incapable of explaining the difference between living organisms and the rest of the natural world. Instead, we need to posit a vital force, or *élan vital*, that explains the distinctive organization, development, and behavior of living things.

Certainly, vitalism has no scientific credibility today. The more that was discovered about the biology and chemistry of living things, the less work there was for an *élan vital*, until finally it became apparent that it was an unnecessary posit because there was no problem to which it might be a solution. But historians of science argue that debates about vitalism served an important role in the development of biology, by forcing biologists to confront some of the explanatory deficiencies of the models they were working with – both by developing new models and by developing new experimental tools. Perhaps mysterianism about the cognitive science of consciousness will have a similar role to play?

Certainly, that would be consistent with how cognitive science has evolved up to now. Many of the advances that we have explored have emerged in response to challenges that on the face of things are no less dramatic than the challenges posed by those who think that consciousness is scientifically inexplicable – the challenge to show how a machine can solve problems, for example; to show how neural networks can learn; or to show how systems can engage in sophisticated emergent behaviors without explicit information processing.

In any event, consciousness is one of the most active and exciting topics in contemporary cognitive science. Whether it will ultimately reveal the limits of cognitive scientific explanation or not, it continues to generate an enormous range of innovative experiments and creative theorizing.



## Summary

This chapter began the basic challenge for cognitive science raised by the apparent conflict first- and third-person approaches to consciousness, as presented in Frank Jackson's knowledge argument. Despite this challenge, consciousness research has made numerous interesting discoveries about the way our minds work. We began with results from priming studies and cases of neurological damage indicating that a great deal of information processing occurs below the threshold of consciousness. Further investigation shows, though, that nonconscious information processing is limited in important ways, and we looked at those limitations to explore the function of consciousness awareness. Milner and Goodale's research on the two visual streams, as well as other related studies, indicate that consciousness is important for planning and initiating actions, while research on planning suggests the memory retains conscious better than nonconscious information. After considering objections to the very idea that cognitive science might provide a complete account of consciousness, we concluded by looking at the global workspace theory of consciousness, which tied together a number of themes in this chapter and throughout the book. The global workspace theory shows how unconscious information reaches consciousness as well as how modular information is transmitted throughout the brain for use in high-level cognition.

## Checklist

### The Challenge of Consciousness

- (1) We can take either a first-person or a third-person approach to consciousness.
- (2) Jackson's Knowledge Argument illustrates the contrast between the first-person and third-person approaches to consciousness.
- (3) According to the Knowledge Argument, a color blind scientist who knew every (third-person) scientific fact about consciousness would still lack a crucial piece of knowledge – first-person knowledge of what it is like to see color
- (4) The contrast between the first- and third-person approaches points to the potential inadequacy of cognitive science for studying consciousness.

### Information Processing without Conscious Awareness

- (1) There are two primary ways of understanding unconscious information processing: priming experiments and studies of patients with neurological damage.
- (2) Semantic priming studies show that basic categorization can be accomplished unconsciously. Since semantic categorization is generally thought to be nonmodular, these tasks also suggest that there can be nonmodular unconscious processing.
- (3) Blindsight and unilateral neglect indicate that high-level processing can be applied even to areas of the visual field that, due to damage, are not consciously perceived.

### The Function of Consciousness

- (1) Milner and Goodale's research reveals a basic functional distinction in the visual system: vision for perception and vision for action. The ventral visual stream is for perception and high-level action-



planning and is conscious, while the dorsal visual stream is for online control of action and is unconscious.

- (2) Experiments on the Ebbinghaus illusion and interocular suppression provide support for Milner and Goodale's dual stream hypothesis.
- (3) Milner and Goodale's research indicates that consciousness is functionally important for planning and initiating action.
- (4) Priming studies show that consciously perceived primes are retained better and have greater impact on other cognitive processes.

### The Hard Problem of Consciousness

- (1) Ned Block's distinction between access consciousness (or A-consciousness) and phenomenal consciousness (or P-consciousness) can be used to generate a dilemma for the cognitive science of consciousness: cognitive science seems to be informative only for understanding A-consciousness.
- (2) Block claims that there is a double dissociation between A- and P-consciousness. This produces an explanatory gap.
- (3) The conflict between A- and P-consciousness can be understood in terms of what David Chalmers calls the hard problem of consciousness, which Chalmers thinks cannot be solved by cognitive science.

### The Global Workspace Theory of Consciousness

- (1) The global workspace theory holds that attention makes low-level modular information available for conscious control in the "global workspace," from where the information is then "broadcast" to other areas of the brain.
- (2) The global workspace theory draws from two basic ideas: (a) consciousness permits information to be explicitly and durably maintained for additional processing and reasoning, and (b) consciousness is necessary for initiating deliberate action.
- (3) Information processing in the global workspace is a type of domain-general process, selecting among competing modular inputs.
- (4) There is some neurological support for the global workspace theory. Pyramidal neurons, for instance, may be responsible for connecting specialized modular processes and broadcasting their outputs throughout the brain for other cognitive processes.

## Further Reading

There has been an explosion of research on consciousness in the last decade or so, only a small portion of which can be covered in a single chapter. Good places to start to learn more are the books by Jesse Prinz 2012 and Timothy Bayne 2012. Though written by philosophers, both books place heavy emphasis on empirical research, and synthesize a wide range of recent studies of consciousness. Robert Van Gulick's chapter in Margolis, Samuels, and Stich 2012 also provides a good summary of both philosophical and neuroscientific theories of consciousness. Zelazo, Moscovitch, and Thompson 2007 is another excellent resource. Baars and Gage 2010 discusses a lot of the most recent research, including figures and descriptions of the most popular methods

used to study consciousness. Chalmers 2013 outlines what he sees as the principal projects for a science of consciousness.

Frank Jackson's Knowledge Argument was first presented in Jackson 1982. His more recent views can be found in Jackson 2003. A series of essays on the Mary thought experiment can be found in Ludlow, Nagasawa, and Stoljar 2004. For a related argument (Joseph Levine's explanatory gap argument), see Levine 1983.

Prominent accounts of how unconscious information processing operates and how information becomes conscious include Dehaene et al. 2006 and Kouider et al. 2007. There are many excellent reviews of research on priming. Kouider and Dehaene 2007 is a good survey of the history of masked priming. On primes becoming more visible when followed by congruent primes, see Bernstein et al. 1989. Good resources on bilingual semantic priming are Kiran and Lebel 2007, Kotz 2001, and Schoonbaert et al. 2009. Classic studies of unilateral neglect include Driver and Mattingly 1998, Driver and Vuilleumier 2001, and Peru et al. 1996. A recent meta-analysis of the critical lesion locations involved in unilateral neglect can be found in Molenberghs, Sale, and Mattingley 2012. Corbetta and Shulman 2011 discusses the relation between neglect and attention. On the function of the parietal cortex in visual perception, see Husain and Nachev 2007.

A summary of the two visual streams can be found in Milner and Goodale 2008. A critique of the two-stream account (with commentary from Milner, Goodale, and others) can be found in Schenk and McIntosh 2010. See Milner 2012 for a recent study of the two visual streams and consciousness. Goodale and Milner 2013 also provides a good review of the visual system. There are many studies on the Ebbinghaus illusion and the differences between vision for action and vision for perception. Aglioti, DeSouza, and Goodale 1995 is a classic study. For responses and follow-up studies, see Glover and Dixon 2001 and Franz et al. 2000.

The literature on access consciousness and phenomenal consciousness is quite large now. Block 1995b is the classic article on the topic. Block's more recent views can be found in Block 2007, where he proposes different neural structures underlying A- and P-consciousness, and Block 2011, where he responds to a number of criticisms of his account. The original Sperling experiment can be found in Sperling 1960. A criticism of Block's interpretation of the Sperling experiment, as well as discussion of phenomenal consciousness more generally, can be found in Kouider et al. 2010. For more on the putative explanatory gap between A- and P-consciousness, see Levine 1983. Other well-known books on these topics include Carruthers 2000 and Dennett 1991. For classic formulations of the hard problem and easy problems of consciousness, see Chalmers 1995, 1996.

For early formulations of the global workspace theory of consciousness, see Baars 1988, 2002. An influential discussion of the theory is Dehaene and Naccache 2001. A summary of the theory can be found in Dehaene and Changeux 2011, including responses to critics. See also Dehaene et al. 2014.

Two popular topics in consciousness research that have been mentioned only briefly but that have their own burgeoning literatures are attention, which was discussed in Chapter 9, and the neural correlates of consciousness. Posner 1980 is a classic early study on attention and consciousness. It was the first to convincingly demonstrate that gaze can be fixed while attention wanders. Lamme 2003 provides a concise summary of the reasons for separating attention from consciousness. Lavie 2005 is an influential account of how unattended stimuli are processed. Mack



and Rock 1998 discusses a series of now-classic experiments on what is called *inattentional blindness*. Simons and Chabris 1999 is another classic series of studies in this area. These experiments rely on selective looking, where people's selective attention alters what they see in a visual array. See Simons and Rensink 2005 for a review of these studies. Other reviews of how attention relates to consciousness can be found in Koch and Tsuchiya 2007, Martens and Wyble 2010, Van den Bussche et al. 2010, and Raffone et al. 2014.

Many trace the most recent wave of research into the neural correlates of consciousness (NCC) to Baars 1988 and Koch 2004. The global workspace theory is one prominent account of the NCC. An influential idea utilized by global workspace theorists is that of neural synchrony. This idea, popularized by Singer 1999, holds that groups of neurons must fire in sync in order to produce consciousness. Womelsdorf et al. 2007 is a more recent paper on this phenomenon. Crick and Koch 2003 is a widely cited review of different problems with the search for NCC, including arguments against the importance of neural synchrony for consciousness. An increasingly popular tool for identifying the NCC is to track brain activation in patients during and after being in a vegetative state. Steven Laureys's studies are some of the best known. Laureys 2005 is an influential article describing the various brain areas that appear to be deactivated as a result of being in a vegetative state. Owen et al. 2006 and Hohwy 2009 are other important articles. Good reviews on the search for the NCC include Metzinger 2000, Lamme 2006, Tononi and Koch 2008, and Koch et al. 2016.





## CHAPTER SIXTEEN

# Robotics: From GOFAI to Situated Cognition and Behavior-Based Robotics

### OVERVIEW 407

#### 16.1 GOFAI Robotics: SHAKEY 408

SHAKY's Software I: Low-Level Activities and Intermediate-Level Actions 409

SHAKY's Software II: Logic Programming in STRIPS and PLANEX 413

#### 16.2 Situated Cognition and Biorobotics 414

The Challenge of Building a Situated Agent 415

#### Situated Cognition and Knowledge

Representation 416

Biorobotics: Insects and Morphological Computation 418

#### 16.3 From Subsumption Architectures to Behavior-Based Robotics 423

Subsumption Architectures: The Example of Allen 424

Behavior-Based Robotics: TOTO 427

Multiagent Programming: The Nerd Herd 430



## Overview

This chapter focuses on robotics. We start by setting the scene with a classic example of GOFAI robotics – GOFAI stands for good old-fashioned artificial intelligence. The focus for the remainder of the chapter is the situated cognition movement in robotics. Like the dynamical systems theorists discussed in Chapter 6, situated cognition theorists are dissatisfied with traditional ways of thinking about information processing in cognitive science. They have developed a powerful tool kit of alternatives and used them to construct new types of robot.

Section 16.1 reviews one of the historic achievements of early robotics, which is also an excellent illustration of the physical symbol system hypothesis in action. This is SHAKEY, a mobile robot developed at the Artificial Intelligence Center at SRI (Stanford Research Institute). SHAKEY was designed to operate and perform simple tasks in a real, physical environment. The programs built into it permitted SHAKEY to plan ahead and to learn how to perform tasks better.

Section 16.2 brings out some of the complaints that situated cognition theorists level at traditional GOFAI robotics, and illustrates some of the engineering inspiration that these theorists have drawn from studying very simple cognitive systems such as insects.

Section 16.3 explores how these theoretical ideas have been translated into particular robotic architectures, focusing on the subsumption architectures developed by Rodney Brooks, and at examples of what Maja Matarić has termed behavior-based robotics.

## 16.1

### GOFAI Robotics: SHAKEY

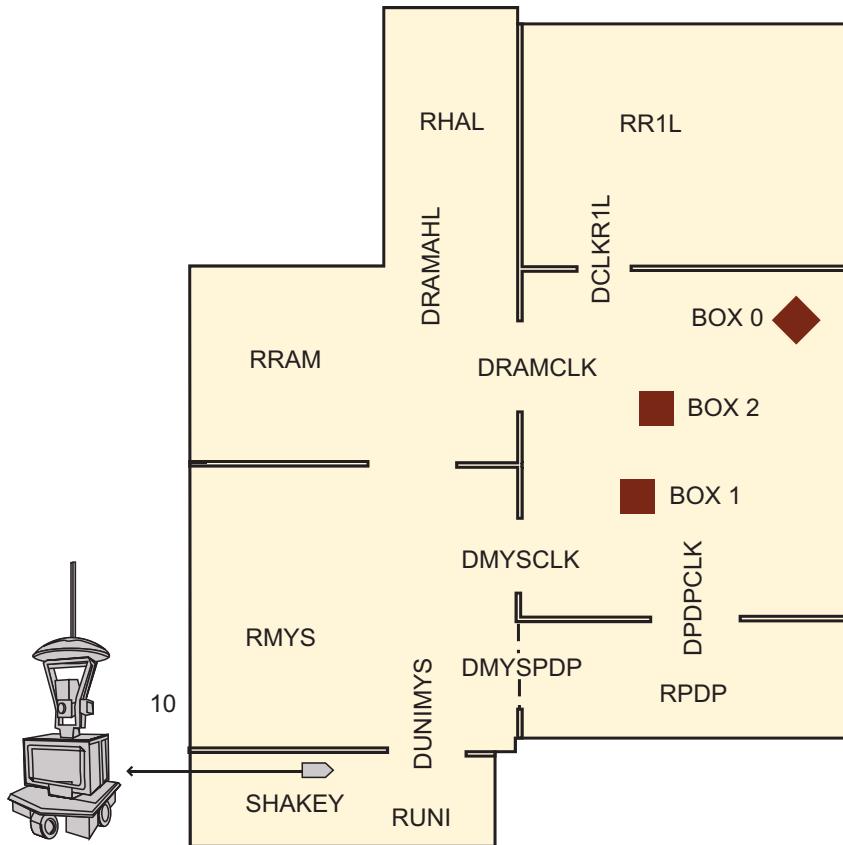
In this section we look at a pioneering robot developed in the late 1960s and early 1970s in the Artificial Intelligence Center at what was then called the Stanford Research Institute (it is now called SRI International and no longer affiliated to Stanford University). This robot, affectionately called SHAKEY (because of its jerky movements), was the first robot able to move around, perceive, follow instructions, and implement complex instructions in a realistic environment (as opposed to virtual micro-worlds like that “inhabited” by SHRDLU, which we looked at in Section 2.1). SHAKEY has now retired from active service and lives in the Robot Hall of Fame at Carnegie Mellon University in Pittsburgh, Pennsylvania.

Figure 16.1 depicts one of the physical environments in which SHAKEY operated. The name of each room begins with an “R.” “RMYS” is a mystery room – i.e., SHAKEY has no information about its contents. Doorway names begin with a “D” and are labeled in a way that makes clear which rooms they are connecting. “DUNIMYS,” for example, labels the door between RUNI (where SHAKEY starts) and RMYS. The environment is empty, except for three boxes located in RCLK (the room with the clock).

In thinking about SHAKEY the first place to start is with the physical structure itself, which is shown in Figure 16.2. The software that allows SHAKEY to operate is not actually run on the robot itself. It was run on a completely separate computer system that communicated by radio with SHAKEY (the radio antenna can be seen in the photo, extending from the top of the structure).

The programs that run SHAKEY are examples of what is generally called *logic programming*. They incorporate a basic model of the environment together with a set of procedures for updating the model and for acting on the environment.

SHAKY’s basic model is given by a set of statements in the *first-order predicate calculus*. (The first-order predicate calculus is the logical language that allows us to talk about particular objects having particular properties, and also permits us to formulate generalizations either about all objects or about at least one object.) These statements are in a basic vocabulary that contains names for the objects in the robot’s world – doors, blocks, walls, and so on – as well as predicates that characterize the properties those objects can have. The vocabulary also contains a name for SHAKEY and predicates that describe the robot’s state – where it is, the angle at which its head is tilted, and so on. The software that SHAKEY uses to plan and execute its actions exploits this same vocabulary, supplemented by terms for particular actions.



**Figure 16.1** A map of SHAKEY's physical environment. Each room has a name. The room containing the boxes is called RCLK (an abbreviation for "Room with Clock"). The total environment measures about 60 feet by 40 feet. (Adapted From Nilsson 1984)

## SHAKEY's Software I: Low-Level Activities and Intermediate-Level Actions

To understand SHAKEY's software we need to go back to Chapter 1, where we looked at Lashley's influential (and at the time very innovative) ideas about the hierarchical organization of behavior. Reacting against the behaviorist idea that actions could be viewed as linked chains of responses, Lashley argued that many complex behaviors resulted from prior planning and organization. These behaviors are organized hierarchically (rather than linearly). An overall plan (say, walking over to the table to pick up the glass) is implemented by simpler plans (the walking plan and the reaching plan), each of which can be broken down into simpler plans, and so on. Ultimately, we arrive at basic actions that don't require any planning. These basic actions are the components from which complex behaviors are built.

SHAKEY's software packages are built around this basic idea that complex behaviors are hierarchically organized. We can see how this works in Table 16.1, which shows how we can think about SHAKEY as a system with five different levels. The bottom level is the hardware level, and there are four different levels of software. The software levels are



**Figure 16.2** A photograph of SHAKEY the robot. (Photograph by Ralph Crane/The LIFE Picture Collection/Getty Images)

hierarchically organized. Each level of software controls a different type of behavior. Going up the hierarchy of software takes us up the hierarchy of behavior.

The interface between the physical hardware of the robot and the software that allows it to act in a systematic and planned way is at the level of low-level actions (LLAs). The LLAs are SHAKEY's basic behaviors – the building blocks from which everything that it does is constructed. The LLAs exploit the robot's basic physical capabilities. So, for example, SHAKEY can move around its environment by rolling forward or backward. It can take photos with the onboard camera and it can move its head in two planes – tilting it up and down, and panning it from side to side. There are LLAs corresponding to all of these abilities.

**TABLE 16.1 SHAKEY'S five levels**

LEVEL	FUNCTION	EXAMPLES
1 Robot vehicle and connections to user programs	To navigate and interact physically with a realistic environment	See the photograph of SHAKEY in Figure 16.2
2 Low-level actions (LLAs)	To give the basic physical capabilities of the robot	ROLL (which tells the robot to move forward by a specified number of feet) and TILT (which tells the robot to tilt its head upward by a specified number of degrees)
3 Intermediate-level actions (ILAs)	Packages of LLAs	PUSH (OBJECT, GOAL, TOL) which instructs the robot to push a particular object to a specified goal, with a specified degree of tolerance
4 STRIPS	A planning mechanism constructing MACROPS (sequences of ILAs) to carry out specific tasks	A typical MACROP might be to fetch a block from an adjacent room
5 PLANEX	Executive program that calls up and monitors individual MACROPS	PLANEX might use the sensors built into the robot to determine that the block can only be fetched if SHAKEY pushes another block out of the way first – and then invoke a MACROP to fetch a block

SHAKY's model of its environment also represents the robot's own state. Of course, executing an LLA changes the robot's state and so requires the model to be updated. Table 16.2 shows the relation between the LLAs that SHAKEY can perform and the way in which it represents its own state.

So, the LLAs fix SHAKEY's basic repertoire of movements. In themselves, however, LLAs are not much use for problem solving and acting. SHAKEY's designers needed to build a bridge between high-level commands (such as the command to fetch a block from a particular room) and the basic movements that SHAKEY can use to carry out that command. The first level of organization above LLAs comes with Intermediate-Level Actions (ILAs). The ILAs are essentially action routines – linked sequences of LLAs that SHAKEY can call upon in order to execute specific jobs, such as navigating to another room, or turning toward a goal. Table 16.3 shows some ILAs.

ILAs are not just chains of LLAs (in the way that behaviorists thought that complex actions are chained sequences of basic responses). They can recruit other ILAs. So, for

**TABLE 16.2** How SHAKY represents its own state

ATOM IN AXIOMATIC MODEL	AFFECTED BY
(AT ROBOT xfeet yfeet)	ROLL
(DAT ROBOT dxfeet dyfeet)	ROLL
(THETA ROBOT degreesleftofy)	TURN
(DTHETA ROBOT dthetadegrees)	TURN
(WHISKERS ROBOT whiskerword)	ROLL, TURN
(OVRID ROBOT overrides)	OVRID
(TILT ROBOT degreesup)	TILT
(DTILT ROBOT ddegreesup)	TILT
(PAN ROBOT degreesleft)	PAN
(DPAN ROBOT ddegreesleft)	PAN
(IRIS ROBOT evs)	IRIS
(DIRIS ROBOT devs)	IRIS
(FOCUS ROBOT feet)	FOCUS
(DFOCUS ROBOT dfeet)	FOCUS
(RANGE ROBOT feet)	RANGE
(TVMODE ROBOT tvmode)	TVMODE
(PICTURESTAKEN ROBOT $\wedge$ picturestaken)	SHOOT

example, the GETTO action routine takes SHAKY to a specific room. This action routine calls upon the NAVTO routine for navigating around in the room SHAKY is currently in, as well as the GOTOROOM routine, which takes SHAKY to the room it is aiming for. Of course, SHAKY can only move from any room to an adjacent room. And so the GOTO-ROOM routine is built up from the GOTOADJROOM routine.

SHAKY's hierarchical organization is very clear even at the level of ILAs. But in order to appreciate it fully we need to look at the next level up. Nothing that we have seen so far counts as planning. Both LLAs and ILAs allow SHAKY to implement fairly low-level commands. But there is little here that would properly be described as problem solving – or indeed, to go back to Newell and Simon's physical symbol system hypothesis, as intelligent action.

**TABLE 16.3** SHAKEY's intermediate-level routines

ILA	ROUTINES CALLED	COMMENTS
PUSH3	PLANOBMOVE*, PUSH2	Can plan and execute a series of PUSH2s
PUSH2	PICLOC*, OBLOC*, NAVTO, ROLLBUMP, PUSH1	Check if object being pushed slips off
PUSH1	ROLL*	Basic push routine; assumes clear path
GETTO	GOTOROOM, NAVTO	Highest-level go-to routine
GOTOROOM	PLANTOUR*, GOTOADJROOM	Can plan and execute a series of GOTOADJROOMs
GOTOADJROOM	DOORPIC*, ALIGN, NAVTO, BUMBLETHRU	Tailored for going through doorways
NAVTO	PLANJOURNEY*, GOTO1	Can plan and execute a trip within one room
GOTO1	CLEARPATH*, PICDETECTOB*, GOTO	Recovers from errors due to unknown objects
GOTO	PICLOC*, POINT, ROLL2	Executes single straight-line trip
POINT	PICTHETA*, TURN2	Orients robot toward goal
TURN2	TURNBACK*, TURN1	Responds to unexpected bumps
TURN1	TURN*	Basic turn routine; expects no bumps
ROLL2	ROLLBACK*, ROLL1	Responds to unexpected bumps
ROLL1	ROLL*	Basic roll routine that expects no bumps
ROLLBUMP	ROLLBACK*, ROLL1	Basic roll routine that expects a terminal bump

## SHAKEY's Software II: Logic Programming in STRIPS and PLANEX

The real innovation in SHAKEY's programming came with the STRIPS planner ("STRIPS" is an acronym for "Stanford Research Institute Problem Solver"). The STRIPS planner (which, as it happens, was fairly closely related to Newell and Simon's General Problem Solver [GPS]) allows SHAKEY to do things that look much more like reasoning about its environment and its own possibilities for action. What STRIPS does is translate a particular goal statement into a sequence of ILAs.

As we have seen, SHAKEY has an axiomatic model of its environment. The axioms are well-formed formulas in the predicate calculus, built up from a basic vocabulary for describing SHAKEY and its environment. These formulas describe both SHAKEY's physical

environment and its own state. The model is updated as SHAKEY moves around and acts upon the environment.

The tasks that SHAKEY is given are also presented as formulas in the predicate calculus. The predicate calculus is a tool for deduction and what SHAKEY does, in essence, is to come up with a deduction that has its goal as its conclusion.

We can think about SHAKEY's planning process as involving a tree search. (Think back to the decision trees that we looked at in Section 12.1.) The first node (the top of the tree) is SHAKEY's model of the current environment. Each branch of the tree is a sequence of ILAs. Each node of the tree is an updated model of the environment.

If the goal formula can be deduced from the updated model at a given node, then STRIPS has solved the problem. What it then does is instruct SHAKEY to follow the sequence of ILAs described in the branch that leads to the solution node. SHAKEY does this, updating its model of the environment as it goes along.

There is no guarantee that this will always get SHAKEY to where it wants to go. The goal might not be attainable. Its model of the environment might not be correct. Someone might have moved the block without telling SHAKEY (and in fact researchers at SRI did do precisely that to see how SHAKEY would update its model). This is where the PLANEX level comes into play. The job of the PLANEX software is to monitor the execution of the plan. So, for example, PLANEX contains an algorithm for calculating the likely degree of error at a certain stage in implementing the task (on the plausible assumption that executing each ILA would introduce a degree of "noise" into SHAKEY's model of the environment). When the likely degree of error reaches a certain threshold, PLANEX instructs SHAKEY to take a photograph to check on its position. If a significant error is discovered, then PLANEX makes corresponding adjustments to the plan.

We can see, then, how SHAKEY is a great example of the physical symbol system hypothesis in action. The physical symbol structures are well-formed formulas in the predicate calculus. These symbols give SHAKEY's model of the environment, as well as the goals and subgoals that SHAKEY is trying to achieve. And the physical symbol structures are manipulated and transformed by the STRIPS and PLANEX algorithms.

Moreover, SHAKEY clearly illustrates the heuristic search hypothesis, also discussed in Chapter 4. The hypothesis says that intelligent problem solving takes place by transforming and manipulating symbol structures until a solution structure is reached. The starting point is given by SHAKEY's model of the environment, together with the target formula that represents the desired end-state. The permissible transformations are given by the STRIPS and PLANEX algorithms. And it is easy to see what the solution structure is. The problem is solved when the initial symbol structure has been transformed into a symbol structure from which the target formula can be deduced.

## 16.2

### Situated Cognition and Biorobotics

With the classic example of SHAKEY clearly in view, we turn now to criticisms of GOFAI robotics and the alternative approaches that have been developed.



We'll start with the situated cognition movement. This is in many ways similar to the dynamical systems approach, which we discussed in Chapter 6. Cognitive scientists influenced by dynamical systems theory propose a new way of analyzing and predicting cognitive systems as coupled systems, focusing on variables evolving through state space in real time, rather than representations. Instead of abstracting away from the physical details of how cognitive systems actually work, they suggest that those physical details can play all sorts of unsuspected but vitally important roles in determining how a cognitive system changes and evolves over time.

Situated cognition theorists propose a broadly similar approach to robotics, moving away from the disembodied approach that we find in SHAKEY (with a sharp distinction between hardware and software) toward design principles that exploit and leverage the robot's physical properties, as well as its embeddedness in real-life environment and its need to operate in real time.

## The Challenge of Building a Situated Agent

The principal objection that situated cognition theorists make to traditional cognitive science is that it has never really come to terms with the real-life problems and challenges in understanding cognition. This problem is particularly acute for the GOFAI approach to building artificial agents.

SHAKY, which we have just looked at, is a classic example of GOFAI robotics. So too is Terry Winograd's SHRDLU program for natural language understanding, which we looked at in Section 2.1. SHRDLU is a virtual robot, reporting on and interacting with a virtual micro-world. A good entry point for the worries that situated cognition theorists have about GOFAI is via a criticism often leveled at SHRDLU and other micro-world programs.

The basic complaint is that SHRDLU only works because its artificial micro-world environment has been stripped of all complexity and challenge. Here is a witty expression of the worry from the philosopher and cognitive scientist John Haugeland (although he is not himself a promoter of the situated cognition movement):

SHRDLU performs so glibly only because his domain has been stripped of anything that could ever require genuine wit or understanding. Neglecting the tangled intricacies of everyday life while pursuing a theory of common sense is not like ignoring friction while pursuing the laws of motion; it's like throwing the baby out with the bathwater. A round frictionless wheel is a good approximation of a real wheel because the deviations are comparatively small and theoretically localized: the blocks-world "approximates" a playroom more as a paper plane approximates a duck.

(Haugeland 1985: 190)

One might wonder whether Haugeland is being completely fair here. After all, Winograd did not really set out to provide "a theory of common sense," and there probably are situations in which a paper plane is a useful approximation of a duck. But the basic point is clear enough. There are many challenges that SHRDLU simply does not have to deal with.

SHRDLU does not have to work out what a block is, for example – or how to recognize one. There is very little "physical" challenge involved in SHRDLU's (virtual) interactions

with its micro-world environment, since SHRDLU has built into it programs for picking up blocks and moving them around, and the robot hand is expressly designed for implementing those programs. Likewise, SHRDLU's language-understanding achievements are partly a function of its artificially limited language and the highly circumscribed conversational context. The major problems in language understanding (such as decoding ambiguity and working out what a speaker is really trying to say) are all factored out of the equation. Finally, SHRDLU is not autonomous – it is a purely reactive system, with everything it does a response to explicit instructions.

In other words, SHRDLU is not properly *situated* in its environment – or rather, the way in which SHRDLU is situated in its environment is so radically different from how we and other real-life cognitive agents are embedded in our environments that we can learn nothing from SHRDLU about how our own cognitive systems work. In fact (the argument continues), SHRDLU's environment is so constrained and devoid of meaning that it is positively misleading to take it as a starting point in thinking about human cognition. The call for situated cognition, then, is a call for AI to work on systems that have all the things that SHRDLU lacks – systems that are properly embodied and have real autonomy. These systems need to be embedded in something much more like the real world, with ambiguous, unpredictable, and highly complex social and physical contexts.

The researchers who designed and built SHAKEY may have thought that they were programming something much closer to an embodied and autonomous agent. After all, SHAKEY can navigate the environment, and it is designed to solve problems, rather than to be purely reactive. But, from the perspective of situated cognition theorists, SHAKEY is really no better than SHRDLU.

For situated cognition theorists, SHAKEY is not really a situated agent, even though it propels itself around a physical environment. The point for them is that the real work has already been done in writing SHAKEY's program. SHAKEY's world is already defined for it in terms of a small number of basic concepts (such as BOX, DOOR, and so forth). Its motor repertoire is built up out of a small number of primitive movements (such as ROLL, TILT, PAN). The problems that SHAKEY is asked to solve are presented in terms of these basic concepts and primitive movements (as when SHAKEY is asked to fetch a BOX).

The robot has to work out a sequence of basic movements that will fulfill the command, but that is not the same as a real agent solving a problem in the real world. SHAKEY already has the basic building blocks for the solution. But working out what the building blocks are is perhaps the most difficult part of real-world problem solving. Like SHRDLU, SHAKEY can only operate successfully in a highly constrained environment. Situated cognition theorists are interested in building agents that will be able to operate successfully even when all those constraints are lifted.

## Situated Cognition and Knowledge Representation

Rodney Brooks is a very influential situated cognition theorist, whose paper “Intelligence without representation” is something of a manifesto for situated cognition theorists. Brooks has an interesting diagnosis of where traditional AI has gone wrong.



Brooks points out that classical AI depends crucially on trimming down the type and number of details that a cognitive system has to represent. Here is his example in his own words:

Consider chairs, for example. While these two characterizations are true

(CAN (SIT-ON PERSON CHAIR))

and

(CAN (STAND-ON PERSON CHAIR))

there is really much more to the concept of a chair. Chairs have some flat (maybe) sitting place, with perhaps a back support. They have a range of possible sizes, and a range of possibilities in shape. They often have some sort of covering material – unless they are made of wood, metal or plastic. They sometimes are soft in particular places. They can come from a range of possible styles. In sum, the concept of what a chair is is hard to characterize simply. There is certainly no AI vision program that can find arbitrary chairs in arbitrary images; they can at best find one particular type of chair in arbitrarily selected images.

(Brooks 1991: 399)

Recognizing and interacting with chairs is a complicated business. But the programmer can remove the complications more or less at a stroke – simply by programming into the system a very narrow characterization of what a chair is. The beauty of doing this is that it can make certain types of chair interactions very simple.

If, to continue with Brooks's example, the system has to solve a problem with a hungry person seated on a chair in a room with a banana just out of reach, then the characterization in the program is just what's required. But of course, if the system solves the problem, then this is largely because it has been given all and only the right sort of information about chairs – and because the problem has been presented in a way that points directly to a solution! Here is Brooks again:

Such problems are never posed to AI systems by showing them a photo of the scene. A person (even a young person) can make the right interpretation of the photo and suggest a plan of action. For AI planning systems, however, the experimenter is required to abstract away most of the details to form a simple description of atomic concepts such as PERSON, CHAIR, and BANANA.

But this abstraction process is the essence of intelligence and the hard part of the problem being solved. Under the current scheme, the abstraction is done by the researchers, leaving little for the AI programs to do but search. A truly intelligent program would study the photograph, perform the abstraction itself, and solve the problem.

(Brooks 1991: 399)

This gives us a much clearer view of what situated cognition is supposed to be all about. It's not just a question of designing robots that interact with their environments. There are plenty of ways of doing this that don't count as situated cognition. The basic idea is to develop AI systems and to build robots that don't have the solutions to problems built into them – AI systems and robots that can learn to perform the basic sensory and motor processes that are a necessary precondition for intelligent problem solving.



## Biorobotics: Insects and Morphological Computation

Situated cognition theorists, like dynamical systems theorists, believe that it pays to start small. Dynamical systems theorists often focus on relatively simple motor and cognitive behaviors, such as infant stepping and the A-not-B error. Cognitive scientists in situated robotics are often inspired by cognitively unsophisticated organisms. Insects are very popular. We can get the flavor from the title of another one of Rodney Brooks's influential articles – "Today the earwig, tomorrow man?"

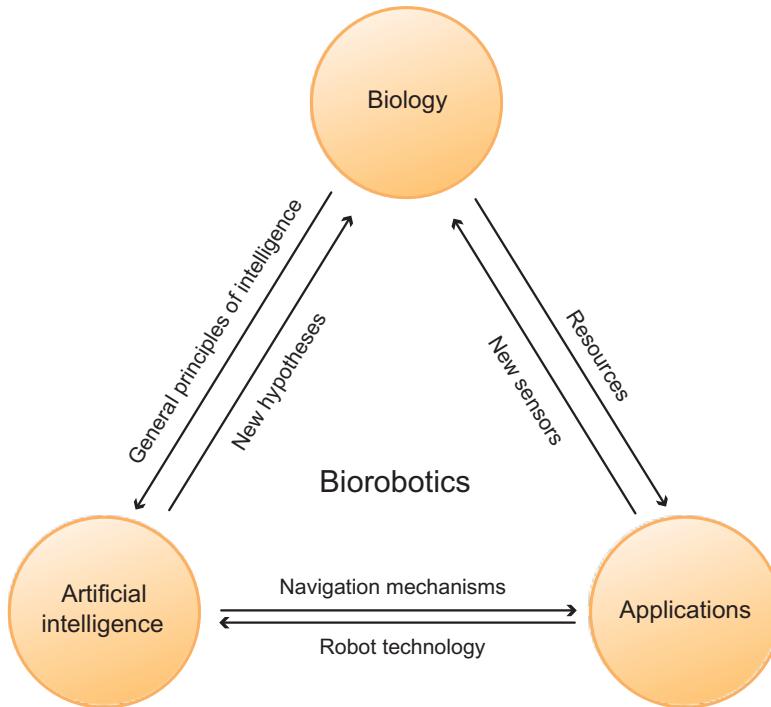
Instead of trying to model highly simplified and scaled-down versions of "high-level" cognitive and motor abilities, situated cognition theorists think that we need to focus on much more basic and *ecologically valid* problems. The key is simplicity without simplification.

Insects solve very complex problems. Studying how they do this, and then building models that exploit the same basic design principles will, according to theorists such as Brooks, pay dividends when it comes to understanding how human beings interact with their environment. We need to look at humans as scaled-up insects, not as scaled-down supercomputers.

One of the basic design principles stressed by situated cognition theorists is that there are direct links between perception and action. This is an alternative to the classical cognitive science view of thinking about organisms in terms of distinct and semi-autonomous subsystems that can be analyzed and modeled independently of each other. On a view like Marr's, for example, the visual system is an autonomous input-output system. It processes information completely independently of what will happen to that information further downstream. When we look at insects, however, we see that they achieve high degrees of "natural intelligence" through clever engineering solutions that exploit direct connections between their sensory receptors and their effector limbs.

Some researchers in this field have described what they are doing as *biorobotics*. The basic idea is usefully summarized in Figure 16.3. Biorobotics is the enterprise of designing and building models of biological organisms that reflect the basic design principles built into those organisms.

Bioroboticists look to biology for insights into how insects and other simple organisms solve adaptive problems, typically to do with locomotion and foraging. They start with theoretical models. They then modify those models after seeing what happens when they are physically implemented in robots – robots whose construction is itself biologically inspired.



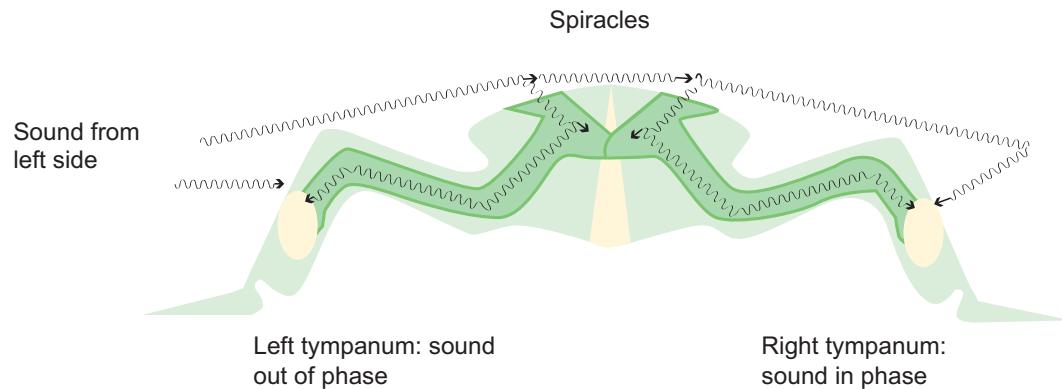
**Figure 16.3** The organizing principles of biorobotics – a highly interdisciplinary enterprise.

A famous example of biorobotics in action is the work of Edinburgh University's Barbara Webb on how female crickets locate males on the basis of their songs – what biologists call cricket phonotaxis.

Female crickets are extremely good at recognizing and locating mates on the basis of the song that they make. On the face of it this might seem a problem that can only be solved with very complex information processing – identifying the sound, working out where it comes from, and then forming motor commands that will take the cricket to the right place. Webb observed, however, that the physiology of the cricket actually provides a very clever solution, exploiting direct links between perception and action.

A remarkable fact about crickets is that they have their ears on their legs. As shown in Figure 16.4, the cricket's ears are connected by a tube (the *tracheal tube*). This means that a single sound can reach each ear via different routes – a direct route (through the ear itself) and various indirect routes (via the other ear, as well as through openings in the tracheal tube known as spiracles). Obviously, a sound that takes the indirect route will take longer to arrive, since it has further to travel – and can't go faster than the speed of sound.

According to Barbara Webb, cricket phonotaxis works because of two very basic design features of the anatomy of the cricket. First, vibration is highest at the ear nearest the source of the sound. This provides a direct indication of the source of the sound. The second is that this vibration directly controls the cricket's movements. Crickets are *hard-wired* to move in the direction of the ear with the highest vibration (provided that the vibration is suitably cricket-like). There is no “direction-calculating mechanism,” no “male cricket identification mechanism,” and no “motor controller.”



**Figure 16.4** The cricket's ears are on its front legs. They are connected to each other via a tracheal tube. The spiracles are small openings that allow air into the tracheal tube. The arrows show the different routes that a single sound can take to each ear. (Adapted from Clark 2001)

Webb and her co-workers have used this model to build robot crickets that can actually perform a version of phonotaxis. In fact, not only can they find the sources of artificial cricket sounds, but they can even find real crickets by their sound.

One of the key design features of Webb's robot cricket (reflecting how real crickets have evolved) is that the cricket's body is a contributing factor in the computation. Cricket phonotaxis works by comparing two different signals. The availability of these two different signals is a direct function of the cricket's bodily layout, as illustrated in Figure 16.4. This can be seen as an early example of what subsequently emerged as the *morphological computation* movement in robotics.

Morphology (in this context) is body shape. The basic idea behind morphological computation is that organisms can exploit features of body shape to simplify what might otherwise be highly complex information-processing tasks. Applying this idea to robotics means building as much of the computation as possible directly into the physical structure of the robot. In essence, morphological computation is a research program for designing robots in which as much computation as possible is done for free.

The morphological computation movement is a very recent development. The first morphological computation conference was only held in 2005. But there have already been some very interesting developments. Here are two examples from the AI Lab in the Department of Informatics at the University of Zurich.

The first example is a fish called WANDA, illustrated in Figure 16.5. WANDA is designed with only one degree of freedom. The only thing WANDA can do is wiggle its tail from side to side at varying amplitudes and frequencies – i.e., WANDA can vary the speed and the degree with which its tail moves. And yet, due to the power of morphological computation, variation in tail wiggling allows WANDA to carry out the full range of fish movements in all three planes – up-down and left-right as well as forward.

Part of the trick here is WANDA's buoyancy, which is set so that slow tail wiggling will make it sink, while fast tail wiggling will make it rise. The other key design feature is the



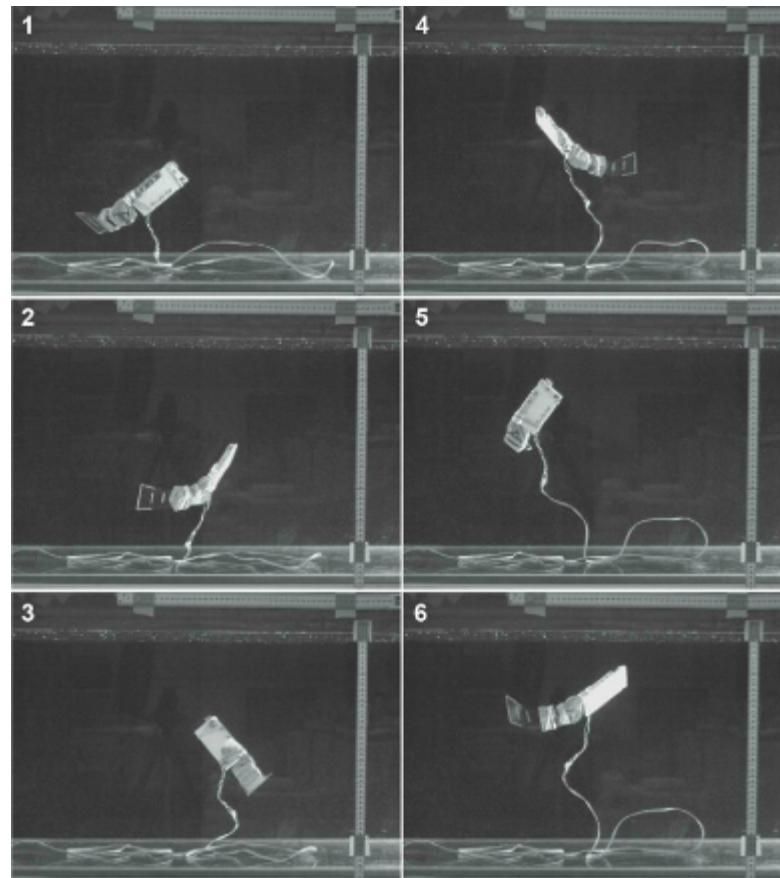
**Figure 16.5** A robot fish called WANDA. All that WANDA can do is wiggle its tail fin. Yet, in an illustration of morphological computation, WANDA is able to swim upward, downward, and from side to side.

possibility of adjusting the zero point of the wiggle movement, which allows for movement to the left or right. Figure 16.6 shows WANDA swimming upward.

A second example of morphological computation also comes from the realm of motor control. (We can think of both examples as ways of counterbalancing the appeal of the computational approach to motor control, which as we saw in Chapter 6 is also a target for dynamical systems theorists) The robot hand devised by Hiroshi Yokoi in Figure 16.7 is designed to avoid the need for making explicit computations in carrying out grasping movements.

On the computational approach, grasping an object requires computing an object's shape and configuring the hand to conform to that shape. Configuring the hand, in turn, requires sending a set of detailed instructions to the tendons and muscles determining the position of the fingers and palm. None of this is necessary, however, in controlling the Yokoi hand.

The hand is constructed from elastic and deformable materials (elastic tendons and deformable fingertips and spaces between the fingers). This morphology does the work that would otherwise be done by complex calculations within some sort of motor control



**Figure 16.6** WANDA swimming upward. (From Pfeifer, Iida, and Gómez 2006)



**Figure 16.7** Another example of morphological computation: the robot hand designed by Hiroshi Yokoi. The hand is partly built from elastic, flexible, and deformable materials. The tendons are elastic, and both the fingertips and the space between the fingers are deformable. This allows the hand to adapt its grasp to the object being grasped.



**Figure 16.8** The Yokoi hand grasping two very different objects. In each case, the control is the same, but the morphology of the hand allows it to adapt to the shapes it encounters. (From Pfeifer, Iida, and Gómez 2006)

unit. What happens is that the hand's flexible and elastic morphology allows it to adapt itself to the shape of the objects being grasped. We see how this works in Figure 16.8.

As with the robot cricket example, most work in morphological computation has focused on the realm of motor control and sensorimotor integration. These are areas where traditional AI, and indeed traditional cognitive science, have often been thought to be deficient. These are not cognitive tasks in any high-level sense. But they are often thought to require information processing, which is why they come into the sphere of cognitive science.

The real question, though, must be how the type of insights that we can find in biorobotics and morphological computation can be integrated into models of more complex agents. Some very suggestive ideas come from the field of behavior-based robotics, to which we turn in the next section.

## 16.3

### From Subsumption Architectures to Behavior-Based Robotics

Rodney Brooks has provided a general AI framework for thinking about some of the agents discussed in the previous section. Webb's robot crickets are examples of what Brooks calls *subsumption architectures*.

Subsumption architectures are organized very differently from the modular architectures that we've been focusing on so far (see Chapter 8, for example) and that are exemplified in SHAKEY. Subsumption architectures are made up of layers and the layers are built up from behaviors.

The bottom level of the architecture is composed of very simple behaviors. Brooks's favorite example is obstacle avoidance, which is obviously very important for mobile robots (and living organisms). The obstacle-avoidance layer directly connects perception (sensing an obstacle) to action (either swerving to avoid the obstacle, or halting where the obstacle is too big to go around).

Whatever other layers are built into the subsumption architecture, the obstacle-avoidance layer is always online and functioning. This illustrates another basic principle

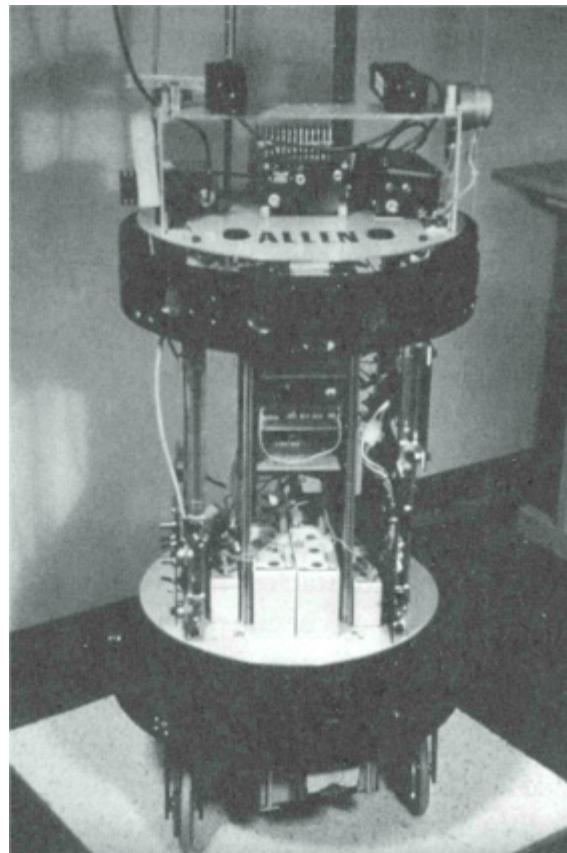
of subsumption architectures. The layers are autonomous and work in parallel. There may be a “higher” layer that, for example, directs the robot toward a food source. But the obstacle-avoidance layer will still come into play whenever the robot finds itself on a collision course with an obstacle. This explains the name “subsumption architecture” – the higher layers subsume the lower layers, but they do not replace or override them.

This makes it easier to design creatures with subsumption architectures. The different layers can be grafted on one by one. Each layer can be exhaustively debugged before another layer is added. And the fact that the layers are autonomous means that there is much less chance that adding a higher layer will introduce unsuspected problems into the lower layers. This is obviously an attractive model for roboticists. It is also explicitly based on thinking about how evolution works.

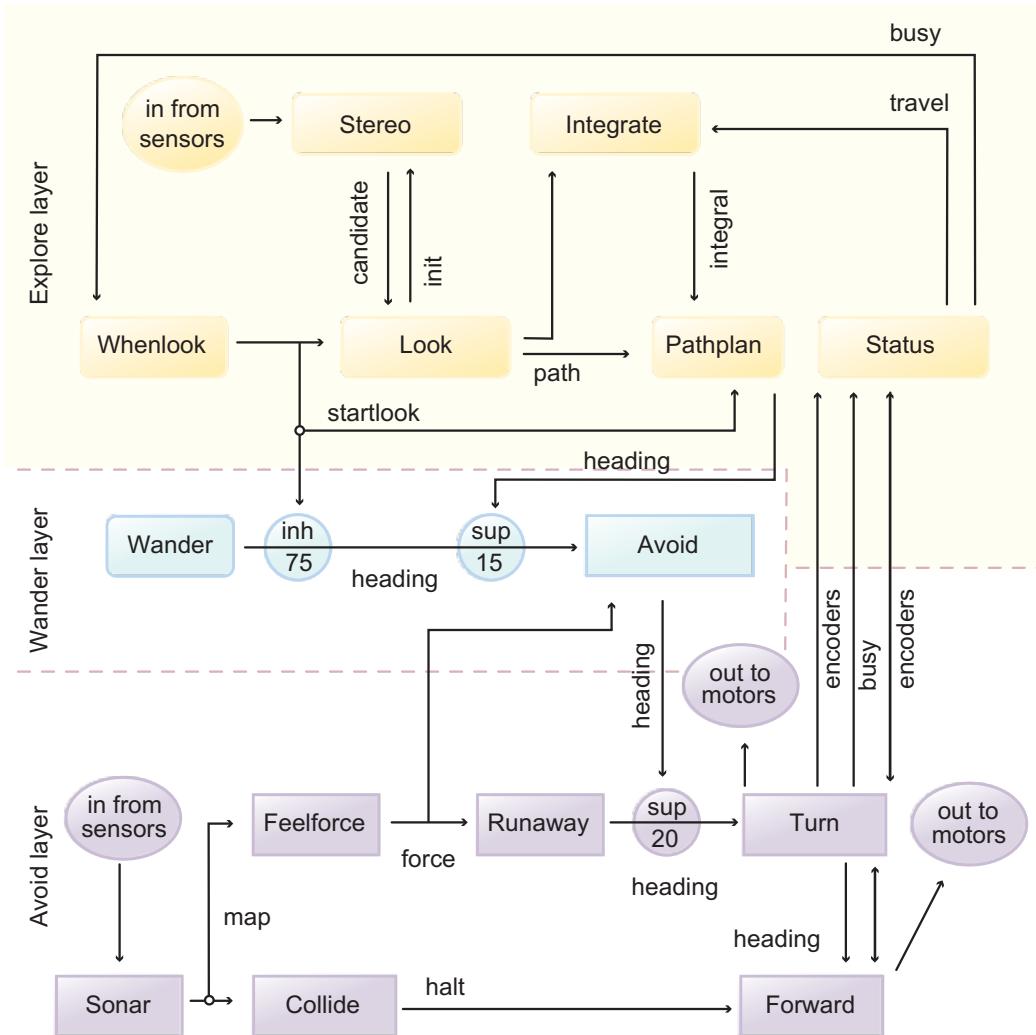
## Subsumption Architectures: The Example of Allen

Rodney Brooks’s lab at MIT has produced many robots with subsumption architectures exemplifying these general principles. One of the first was Allen, illustrated in Figure 16.9.

At the hardware level, Allen looks like any GOFAI robot, not very dissimilar to SHAKEY for example. At the software level, though, Allen is a subsumption architecture, built up in



**Figure 16.9** Rodney Brooks’s robot Allen, his first subsumption architecture robot. (From Brooks 1991)



**Figure 16.10** The layers of Allen's subsumption architecture. Allen has a three-layer architecture. The layers communicate through mechanisms of inhibition (inh) and suppression (sup). (From Brooks 1991)

the standard layered manner. Over time, more and more layers were added to Allen's basic architecture. The first three layers are depicted in Figure 16.10.

The most basic layer is the obstacle-avoidance layer. The diagram shows that this layer is itself built up from a number of distinct subsystems. These do pretty much what their names suggest. The COLLIDE subsystem scans the sensory input for obstacles. It sends out a halt signal if it detects one. At the same time the FEELFORCE system works out the overall force acting upon the robot (using information from the sensors and the assumption that objects function as repulsive forces). These feed into systems responsible for steering the robot – systems that are directly connected to the motor effectors.

The wander and explorer layers are constructed in the same way. In the middle layer the WANDER component generates random paths for Allen's motor system, while the AVOID component feeds back down into the obstacle-avoidance layer to ensure that following the

random path does not lead Allen to crash into anything. Allen is actually pretty successful at this. The robot can successfully navigate environments with both stationary obstacles and other moving objects.

But Allen is not just wandering around at random. The subsystems in the top layer (the explorer layer) work together to allow Allen to pursue goals in a self-directed way. These subsystems receive input from the sensory systems and allow Allen to plan routes toward specific visible locations. As the wiring diagram in Figure 16.10 shows, the PATHPLAN subsystem feeds into the AVOID subsystem. This allows for the plan to be modified as the robot is actually moving toward the goal.

Drawing all this together, we can identify three basic features of subsumption architectures, as developed by Brooks and other AI researchers:

- *Incremental design*: Subsumption architecture robots are built to mimic how evolution seems to work. New subsystems are grafted on in layers that typically don't change the design of the existing subsystems.
- *Semi-autonomous subsystems*: The subsystems operate relatively independently of each other, although some subsystems are set up to override others. The connections between the subsystems are hard-wired. There is typically no central "controller," comparable to STRIPS and PLANEX in SHAKEY
- *Direct perception-action links*: Subsumption architectures trade as much as possible on subsystems that deliver immediate motor responses to sensory input. They are designed for real-time control of action.

The contrast with traditional AI approaches is sharp. Traditional AI robots such as SHAKEY are designed in a very top-down way. There is typically a central planner maintaining a continuously updated model of the world, updated by incorporating information received through its sensors. The planner uses this model of the world to work out detailed action plans, which are transmitted to the effectors. The action plans tend to be multistage and leave little scope for modification.

Proponents of GOFAI robotics are likely to say that the basic features of subsumption architectures are very good design principles for robots that are intended to be no more than mechanical insects – basically capable only of moving around the environment and reacting in simple ways to incoming stimuli. But subsumption architectures are not going to help us with complex intelligent behavior.

Recall the phrasing of the physical symbol system hypothesis, which we looked at in detail in Chapter 4 and which is the dominant theoretical framework for GOFAI robotics. The physical symbol system hypothesis is a hypothesis about the necessary and sufficient conditions of *intelligent action*. But how intelligent is Allen, or the robot crickets and cockroaches that bioroboticists have developed?

Subsumption architectures certainly don't seem to have any decision-making processes built into them. Potential conflicts between different layers and between individual subsystems within a layer are resolved by precedence relations that are built into the hardware



of the robot. Conflict resolution is purely mechanical. But what makes a system intelligent, one might reasonably think, is that it can deal with conflicts that cannot be resolved by applying independent subsystems in some predetermined order. Subsumption architectures lack intelligence almost by definition.

There are different ways in which a situated cognition theorist might try to respond to this challenge. One way is to try to combine the two approaches. There are hybrid architectures that have a subsumption architecture for low-level reactive control, in combination with a more traditional central planner for high-level decision-making. So, for example, Jonathan Connell, a researcher at IBM's T. J. Watson Research Center in Yorktown Heights, New York, has developed a three-level hybrid architecture that he calls SSS.

It is easy to see where the acronym comes from, when we look at what each of the layers does. SSS contains

- a **Servo**-based layer that controls the robot's effectors and processes raw sensory data
- a **Subsumption** layer that reacts to processed sensory input by configuring the servo-based layer (as is standard in a subsumption architecture, the different subsystems are organized in a strict precedence hierarchy)
- a **Symbolic** layer that maintains complex maps of the environment and is capable of formulating plans; the symbolic layer configures the subsumption layer

The hybrid architecture approach abandons some of the basic ideas behind situated cognition and biorobotics. To return to a phrase used earlier, situated cognition theorists like to think of sophisticated cognitive systems as scaled-up insects, whereas GOFAI theorists think of them as scaled-down supercomputers. The hybrid architecture approach, as its name suggests, looks for a middle way – it sets out to build scaled-up insects with scaled-down supercomputers grafted onto them.

But some situated cognition theorists have tried to meet the challenge without compromising on the basic principles of situated cognition. *Behavior-based robotics* moves beyond basic subsumption architectures in a way that tries to build on the basic insights of the situated cognition movement.

## Behavior-Based Robotics: TOTO

Behavior-based architectures are designed to be capable of representing the environment and planning complex actions.

Subsumption architectures (and insects, for that matter) are purely reactive – they are designed to respond quickly to what is happening around them. These responses are typically fairly simple – such as changing the robot's direction, or putting it into reverse when a collision is anticipated. These responses tend to be explicitly programmed in the system.

Behavior-based robots, in contrast, are capable of more complex behaviors that need not be explicitly specified within the system. These are sometimes called *emergent* behaviors

(because they emerge from the operation and interaction of lower-level behaviors). The important thing is that this additional cognitive sophistication is gained without a central planner that works symbolically.

Behavior-based architectures incorporate some of the basic design features of subsumption architectures. They are typically built up from semi-autonomous subsystems in a way that mimics the incremental approach that evolution seems to take. They have two key features, one that they share with subsumption architectures and one that sets them apart:

- *Real-time functioning*: Like subsumption architectures, behavior-based architectures are designed to operate in real time. That is, they make plans on a timescale that interfaces directly with the robot's movements through the environment. This contrasts with symbolic planners and hybrid architectures, where planning is done offline and then needs to be integrated with the robot's ongoing behavior.
- *Distributed representations*: Behavior-based architectures represent their environments and use those representations in planning actions. This distinguishes them from most subsumption architectures. But, unlike symbolic and hybrid architectures, those representations are not centralized or centrally manipulated. There is no central planning system that gathers together all the information that the robot has at its disposal.

We can appreciate how these features work by looking at two examples from the work of Maja Matarić, one of the pioneers of behavior-based robotics. One of the very interesting features of Matarić's work is how she applies the behavior-based approach to programming collections of robots. We will look in some detail at an example of multiagent programming. First, though, let's look briefly at how behavior-based robotics works for single robots.

A fundamental design feature of behavior-based architectures is the distinction between *reactive rules* and *behaviors*. Subsumption architectures are basically built up from reactive rules. A reactive rule might, for example, tell the robot to go into reverse when its sensors detect a looming object. The reactive rules exploit direct perception-action links. They take inputs from the robot's sensors and immediately send instructions to the robot's effectors.

Behaviors, in contrast, are more complex. Matarić defines a behavior as *a control law that satisfies a set of constraints to achieve and maintain a particular goal*. The relevant constraints come both from the sensed environment (which might include other robots) and from the robot itself (e.g., its motor abilities).

So, the challenge for behavior-based robotics is to find a way of implementing behaviors in a mobile agent without incorporating a symbolic, central planner. Matarić's robot TOTO, which she designed and constructed together with Rodney Brooks, illustrates how this challenge can be met for a very specific navigation behavior. This is the behavior of finding the shortest route between two points in a given environment.

Matarić and Brooks were inspired by the abilities of insects such as bees to identify shortcuts between feeding sites. When bees travel from their hive they are typically capable of flying directly to a known feeding site without retracing their steps. In some sense they (and many other insects, foraging animals, and migrating birds) are constructing and



updating maps of their environment. This is a classic example of an apparently complex and sophisticated behavior being performed by creatures with very limited computational power at their disposal – exactly the sort of thing that behavior-based robotics is intended to model.

TOTO is designed to explore and map its environment (an office-like environment where the principal landmarks are walls and corridors) so that it can plan and execute short and efficient paths to previously visited landmarks. TOTO has a three-layer architecture. The first layer comprises a set of reactive rules. These reactive rules allow it to navigate effectively and without collisions in its environment. The second layer (the landmark-detector layer) allows TOTO to identify different types of landmark. In the third layer, information about landmarks is used to construct a distributed map of the environment. This map is topological, rather than metric. It simply contains information as to whether or not two landmarks are connected – but not as to how far apart they are. TOTO uses the topological map to work out in real time the shortest path back to a previously visited landmark (i.e., the path that goes via the smallest number of landmarks).

One of TOTO's key features is that its map is distributed (in line with the emphasis within behavior-based robotics on distributed representations) and the processing works in parallel. There is no single data structure representing the environment. Instead, each landmark is represented by a procedure that categorizes the landmark and fixes its compass direction.

The landmark procedures are all linked together to form a network. Each node in the network corresponds to a particular landmark, and if there is a direct path between two landmarks then there is an edge connecting them in the network. This network is TOTO's topological map of the environment. It is distributed because it exists only in the form of connections between separate landmark procedures.

Behavior-based roboticists do not object to representations per se. They recognize that any robot capable of acting in complex ways in a complex environment must have some way of storing and processing information about its environment. Their real objection is to the idea that this information is stored centrally and processed symbolically.

TOTO constantly expands and updates its network as it moves through the environment detecting new landmarks. This updating is done by activation spreading through the network (not dissimilar to a connectionist network). When the robot is at a particular landmark the node corresponding to that landmark is active. It inhibits the other nodes in the network (which is basically what allows TOTO to know where it is), at the same time as spreading positive activation (expectation) to the next node in the direction of travel (which allows TOTO to work out where it is going).

This distributed map of the environment leaves out a lot of important information (about distances, for example). But for that very reason it is flexible, robust, and, most importantly, very quick to update.

Matarić and Brooks designed an algorithm for TOTO to work out the shortest path between two nodes on the distributed map. The algorithm works by spreading activation.

Basically, the active node (which is TOTO's current location) sends a call signal to the node representing the target landmark. This call signal gets transmitted systematically through the network until it arrives at the target node. The algorithm is designed so that the route that the call signal takes through the network represents the shortest path between the two landmarks. Then TOTO implements the landmark procedures lying on the route to navigate to the target landmark.

TOTO is a nice example of the key features of behavior-based robotics. TOTO is not simply a reactive agent, like Barbara Webb's robot cricket. Nor does it have a central symbolic planner like Jonathan Connell's SSS. It is capable of fairly sophisticated navigation behavior because it has a distributed map of the environment that can be directly exploited to solve navigational problems. The basic activation-spreading mechanisms used for creating and updating the map are the very same mechanisms used for identifying the shortest paths between two landmarks. The mechanisms are somewhat rough-and-ready. But that is what allows them to be used efficiently in the real-time control of behavior – which, after all, is what situated cognition is all about.

## Multiagent Programming: The Nerd Herd

Multiagent programming is highly demanding computationally, particularly if it incorporates some sort of centralized planner or controller. A central planner would need to keep track of all the individual robots, constantly updating the instructions to each one to reflect the movements of others – as well as the evolution of each robot's own map of the environment. The number of degrees of freedom is huge.

The multiagent case presents in a very stark way the fundamental challenges of robotics. How can one design a system that can reason about its environment without a complete combinatorial explosion? It is very instructive to see what happens when the challenge is tackled through the behavior-based approach.

Matarić built a family of twenty mobile robots – the so-called Nerd Herd, illustrated in Figure 16.11. Each robot was programmed with a set of basis behaviors. These basis



**Figure 16.11** The Nerd Herd, together with the pucks that they can pick up with their grippers.

**TABLE 16.4** The five basis behaviors programmed into Matarić's Nerd Herd robots

<b>Safe-wandering</b>	Ability to move around while avoiding collisions with robots and other objects
<b>Following</b>	Ability to move behind another robot retracing its path
<b>Dispersion</b>	Ability to maintain a <i>minimum</i> distance from other robots
<b>Aggregation</b>	Ability to maintain a <i>maximum</i> distance from other robots
<b>Homing</b>	Ability to find a particular region or location

behaviors served as the building blocks for more complex *emergent* behaviors that were not explicitly programmed into the robots.

Table 16.4 shows the five basis behaviors that Matarić programmed into the robots in the Nerd Herd. These behaviors could be combined in two ways. The first way is through summation. The outputs from two or more behaviors are summed together and channeled toward the relevant effector (e.g., the wheels of the robot). This works because all of the behaviors have the same type of output. They all generate velocity vectors, which can easily be manipulated mathematically. The second combination is through switching. Switching inhibits all of the behaviors except for one.

Each of these basis behaviors is programmed at the level of the individual robot. None of the basis behaviors is defined for more than one robot at a time and there is no communication between robots. What Matarić found, however, was that combining the basis behaviors at the level of the individual robots resulted in emergent behaviors at the level of the group.

So, for example, the Nerd Herd could be made to display flocking behavior by summing basis behaviors in each individual robot. The group flocked together as a whole if each robot's control architecture summed the basis behaviors Dispersion, Aggregation, and Safe-wandering. Adding in Homing allowed the flock to move together toward a particular goal.

The principal activity of the robots in the Nerd Herd is collecting little pucks. Each robot has grippers that allow it to pick the pucks up. Matarić used the control technique of switching between different basis behaviors in order to generate the complex behavior of foraging. If the robot doesn't have a puck then all the basis behaviors are inhibited except Safe-wandering. If Safe-wandering brings it too close to other robots (and hence to potential competitors) then the dominant behavior switches to Dispersion. If it has a puck then the control system switches over to Homing and the robot returns to base.

You may be wondering just how intelligent these complex behaviors really are. It is true that flocking and foraging are not explicitly programmed into the system. They are emergent in the sense that they arise from the interaction of basis behaviors. But the mechanisms of this interaction are themselves programmed into the individual robots using the combinatorial operators for basis behaviors. They are certainly not emergent in

the sense of being unpredictable. And one might think that at least one index of intelligence in robots or computers more generally is being able to produce behaviors that cannot simply be predicted from the wiring diagram.

It is significant, therefore, that Matarić's behavior-based robots are capable of learning some of these complex behaviors without having them explicitly programmed. She showed this with a group of four robots very similar to those in the Nerd Herd. The learning paradigm she used was reinforcement learning. What are reinforced are the connections between the states a robot is in and actions it takes.

The complex behavior of foraging is really just a set of condition–behavior pairs – if the robot is in a certain condition (e.g., lacking a puck) then it yields total control to a single behavior (e.g., Safe-wandering). So, learning to forage is, in essence, learning these condition–behavior pairs. This type of learning can be facilitated by giving the robot a reward when it behaves appropriately in a given condition, thus reinforcing the connection between condition and behavior.

Matarić worked with two types of reinforcement – reinforcement at the completion of a successful behavior, and feedback while the robot is actually executing the behavior. Despite the complexity of the environment and the ongoing multiagent interactions, Matarić found that her four robots successfully learned group foraging strategies in 95 percent of the trials.

Obviously, these are early days for behavior-based robotics. It is a long way from groups of robots foraging for hockey pucks in a closed environment to anything recognizable as a human social interaction. But behavior-based robotics does at least give us a concrete example of how some of the basic insights behind the situated cognition movement can be carried forward. Perhaps it is time to change Rodney Brooks's famous slogan: "Yesterday the earwig. Today the foraging robot. Tomorrow man?"



## Summary

After reviewing the classic example of SHAKEY in GOFAI robotics, this chapter focused on new developments in designing and building artificial agents. After reviewing some of the objections that situated cognition theorists level at traditional GOFAI we explored how these theorists have been inspired by very simple cognitive systems such as insects. We then considered how these theoretical ideas have been translated into particular robotic architectures, focusing on the subsumption architectures developed by Rodney Brooks and on Maja Matarić's behavior-based robotics.

## Checklist

**The robot SHAKEY is an example of how a physical symbol system can interact with a real physical environment and reason about how to solve problems.**



- (1) SHAKEY has a model of its environment given by a set of sentences in a first-order logical language. This model is updated as SHAKEY moves around.
- (2) SHAKEY's software is hierarchically organized into four different levels. At the most basic level are primitive actions (low-level actions – LLAs). These LLAs are organized into action routines (Intermediate-level actions – ILAs). SHAKEY solves problems by constructing a sequence of ILAs that will achieve a specific goal.
- (3) The STRIPS problem-solving software is an example of logic programming. It explores the logical consequences of SHAKEY's model of its current environment in order to work out which ILAs can be applied in that environment.
- (4) STRIPS then works out how the model would need to be updated if each ILA were executed in order to develop a tree of possible ILA sequences. If one of the branches of the tree leads to the desired goal state then SHAKEY implements the sequence of ILAs on that branch.

**Situated cognition theorists react against some of the fundamental tenets of GOFAI cognitive science and robotics.**

- (1) GOFAI programs such as SHRDLU and SHAKEY can interact (virtually) with their environments. But situated cognition theorists argue that they are not properly situated in their environments. The real work of putting intelligence into these agents is not done by the systems themselves, but by the programmers.
- (2) The world of a GOFAI robot is already defined for it in terms of a small number of basic concepts. Likewise for its motor repertoire. This avoids the real problems of decoding the environment and reacting to the challenges it poses.
- (3) Situated cognition theorists think that instead of focusing on simplified and scaled-down versions of "high-level" tasks, cognitive scientists should look at how simple organisms such as insects solve complex but ecologically valid problems.
- (4) Biorobotics is the branch of robotics that builds models of biological organisms reflecting the basic design principles that have emerged in evolution. A good example is Barbara Webb's work on cricket phonotaxis.

**Subsumption architectures are a powerful tool developed by situated cognition theorists such as Rodney Brooks.**

- (1) Subsumption architectures are not made up of functional subsystems in the way that modular architectures are. Instead they are built up from layers of semi-autonomous subsystems that work in parallel.
- (2) Subsumption architectures are built to mimic how evolution might work. New systems are grafted on in layers that typically don't change the design of the existing systems.
- (3) Subsumption architectures trade as much as possible on direct perception-action links that allow the online control of action.

**Subsumption architectures do not typically have decision-making systems built into them. Problems of action selection are solved by predefined precedence relations among**

**subsystems. Situated cognition theorists have to work out a more flexible solution to the action selection problem**

- (1) One approach is to develop a hybrid architecture, combining a subsumption architecture for low-level reactive control with a more traditional symbolic central planner for high-level decision-making.
- (2) Behavior-based robotics takes another approach, more in the spirit of situated cognition. Behavior-based architectures (such as that implemented in TOTO) represent their environments and use those representations to plan actions. But these representations are not centralized or centrally manipulated.
- (3) In addition to reactive rules such as those in subsumption architectures, behavior-based robots have basic behaviors programmed into them. These basic behaviors are more complex and temporally extended than reactive rules. They can also be combined.
- (4) Behavior-based robots can exhibit emergent behaviors that have not been programmed into them (e.g., the flocking and foraging behaviors displayed by Matarić's Nerd Herd). Behavior-based robots have also been shown to be capable of learning these emergent behaviors through reinforcement learning.

## Further Reading

SHAKY is very well documented in technical reports published by SRI. These can be downloaded at [www.ai.sri.com/shakey/](http://www.ai.sri.com/shakey/). Technical report 323 is particularly helpful. Also see the *Encyclopedia of Cognitive Science* entry on STRIPS.

The philosopher Andy Clark is a very clear expositor of situated cognition and biorobotics – see particularly his book *Being There* (1997) and chapter 6 of Clark 2001, as well as his book *Supersizing the Mind* (2008) and a discussion of the book in *Philosophical Studies* (2011). For more on morphological computation, including the two examples discussed in the text, see Pfeifer, Iida, and Gómez 2006. Clancey 1997 is a general survey of situated cognition from the perspective of an AI specialist. Several of Rodney Brooks's influential papers are reprinted in his book *Cambrian Intelligence* (1999), which also contains some more technical papers on specific architectures. Brooks 1991 is also reprinted in Haugeland 1997.

For early versions of some of the criticisms of GOFAI made by situated cognition theorists, see Dreyfus 1977. For a very different way of thinking about situated cognition (in terms of situatedness within a social environment), see Hutchins 1995. The *Cambridge Handbook of Situated Cognition* (Robbins and Aydede 2008) is a useful and comprehensive resource, with a strong emphasis on the philosophical underpinnings of the situated cognition movement. For more on embodied cognition, see Shapiro 2007, Chemero 2009, Adams and Aizawa 2010, Shapiro 2011, Anderson, Richardson, and Chemero 2012, and Lawrence Shapiro's chapter in Margolis, Samuels, and Stich 2012.

Arkin 1998 is a comprehensive textbook on behavior-based robotics. For a more programming-oriented survey, see Jones and Roth 2003. Winfield 2012 is a more recent introduction. Maja Matarić has written many papers on behavior-based robotics (see online resources). Matarić 1997,



1998 are good places to start. Readers interested in building their own mobile robots will want to look at her book *The Robotics Primer* (2007).

The robot cricket made its first appearance in Webb 1995. Shigeo Hirose was a pioneer of bio-inspired robotics, starting out with snakes, as described in his book Hirose 1993. For a more contemporary overview of bio-inspired robotics, see Ijspeert 2014 and, for examples of specific projects, see Libby et al. 2012 (tailed-assisted pitch control in lizards and dinosaurs). Koh et al. 2015 (water jumping robots) and Graule et al. 2016 (insects perching and taking off).





## CHAPTER SEVENTEEN

# Looking Ahead: Challenges and Opportunities

### OVERVIEW 437

- 17.1 Exploring the Connectivity of the Brain: The Human Connectome Project and Beyond** 438
- 17.2 Understanding What the Brain Is Doing When It Appears Not to Be Doing Anything** 439

- 17.3 Neural Prosthetics** 440
- 17.4 Cognitive Science and the Law** 441
- 17.5 Autonomous Vehicles: Combining Deep Learning and Intuitive Knowledge** 442



## Overview

Cognitive science has provided massive and important insights into the human mind. We have explored a good number of these in this book. These insights all stem from the single basic idea governing cognitive science as the interdisciplinary science of the mind. This is the idea that mental operations are information-processing operations.

This book began by looking at how this way of thinking about the mind first emerged out of developments in seemingly disparate subjects, such as mathematical logic, linguistics, psychology, and information theory. Most of the significant early developments in cognitive science explored the parallel between information processing in the mind and information processing in a digital computer. As cognitive scientists and cognitive neuroscientists developed more sophisticated tools for studying and modeling the brain, the information-processing principle was extended in new directions and applied in new ways.

The interdisciplinary enterprise of cognitive science is now in excellent health. There are more contributing disciplines than ever before. Cognitive scientists have an ever-expanding range of theoretical models to work with. And there is a constant stream of technological advances in the machinery that cognitive scientists can use to study the brain. It is hard not to have a sense of optimism – a sense that cognitive science is getting close to a fundamental breakthrough in understanding cognition and the mind.



What I want to do now is to look ahead at some of the challenges and opportunities facing cognitive science at this exciting time. What follows is a small and highly personal selection of these challenges and potential applications.

## Exploring the Connectivity of the Brain: The Human Connectome Project and Beyond

The successful completion of the Human Genome Project was one of the most significant scientific events of the last few decades. For the first time scientists succeeded in identifying and mapping the 20,000 to 25,000 genes in the human gene pool, giving unprecedented insights into human genetic makeup. The Human Genome Project was so successful that it focused the minds of funding agencies on huge, collaborative projects.

In July 2009 the National Institutes of Health (NIH) announced what is in effect a cognitive science equivalent of the Human Genome Project – the *Human Connectome Project*. According to the funding opportunity announcement, “The overall purpose of this 5-year Human Connectome Project (HCP) is to develop and share knowledge about the structural and functional connectivity of the human brain.” Awards of \$40 million were made to two consortia of research universities and institutes. The resulting collaborative and multisite effort has tackled some of the theoretical issues highlighted at various points in this book – such as the relation between different types of brain connectivity, and the importance of calibrating different tools for studying the brain.

As outlined in a position paper in *Nature Neuroscience* by leading researchers in one of the two funded consortia (Glasser et al. 2016), the HCP has generated not just considerable amounts of data, but also proposed a new paradigm for neuroimaging research that aims to address longstanding methodological problems with scanning and analyzing the resulting data (some of which were identified in Chapter 9). Key elements of the paradigm proposed in Glasser et al. 2016 include

- acquiring large amounts of high-quality data on as many subjects as feasible, combining different experimental techniques
- focusing on data with high spatial and temporal resolution, and removing distortions, noise and temporal artifacts
- representing cortical and subcortical neuroimaging data in a common geometrical framework (*brainordinates*), represented in a distinctive file format (CIFTI)
- developing a *parcellation* of the brain into distinct regions, based on connectivity and neuroanatomy
- routinely sharing extensively analyzed results such as statistical maps (plus raw and preprocessed data when feasible) together with the code used for the analysis, so that other neuroscientists can make precise comparisons across studies, along with replicating and extending findings

The aim and promise is to make neuroimaging more standardized and systematic. This will make comparisons across subjects and populations much easier. In particular, it will make it much easier to compare normal brains with brains that have suffered damage or disease.



Looking ahead, important challenges include incorporating this (or a comparable) standardizing framework across neuroimaging research, and then, within that framework, developing databases on specific disorders. At the time of writing (summer 2018), ongoing projects under the auspices of the Connectome Coordination Facility include the Alzheimer's Disease Connectome Project, the Epilepsy Connectome Project, the Human Connectome Project for Early Psychosis, and Connectomes Related to Anxiety and Depression.

## 17.2

# Understanding What the Brain Is Doing When It Appears Not to Be Doing Anything

One of the themes of the Human Connectome Project has been the importance of studying the brain's resting state. This is another frontier for neuroscience.

Neuroimaging and electrophysiological experiments standardly explore what happens in the brain when certain very specific tasks are being carried out. So, for example, neuroimaging experiments typically identify the different brain areas where the BOLD contrast is highest during a given task. This is the basis for inferences about localization of function in the brain. But, some researchers have argued, task-specific activation is simply the tip of the iceberg. Marcus Raichle and colleagues at Washington University in St. Louis have argued that we shouldn't just pay attention to departures from the baseline set by the brain's default mode of operation. There is a huge amount of activity going on in the brain even when subjects are resting with their eyes closed, or passively looking at a fixed stimulus. This default mode of brain function has not yet been systematically studied by neuroscientists, but may be quite fundamental to understanding cognition. Concentrating solely on task-dependent changes in the BOLD signal may turn out to be like trying to understand how tides work by looking at the shape of waves breaking on the shore.

What is now often called the *default mode network* (DMN) can be studied in pure resting state experiments, where subjects are imaged while not performing any directed task. The brain areas most frequently identified in such experiments include the medial posterior cortex, particularly the posterior cingulate cortex and the precuneus, and the medial frontal cortex, in addition to areas around the temporoparietal junction area (TPJ).

One interesting possibility starting to gain traction is that some cognitive disorders and diseases may be correlated with impaired functioning of the DMN. A number of studies have observed significant correlations between deteriorating connectivity of the DMN over time and two well-known markers of early Alzheimer's – rising levels of amyloid beta (the key component of brain plaques in Alzheimer's) and pathologies of tau protein (which form tangles inside neurons and disturbs synaptic communication). For large studies, see Chhatwal et al. 2013, Thomas et al. 2014, and Buckley et al. 2017. Epilepsy, Parkinson's, and ADHD are other disorders where impaired functioning of the DMN may be important, as reviewed in Mohan et al. 2016.

 17.3

## Neural Prosthetics

Suppose that, as many cognitive scientists think, important cognitive functions are carried out by functionally specialized systems that are themselves implemented in specific neural locations. Wouldn't it then be possible to build mechanical devices that could replace a damaged system in the brain, reversing the effects of disease or injury? Cognitive science certainly predicts that this type of *neuroprosthesis* ought to be possible. If cognitive systems are computational devices, whose job is basically transforming a certain type of input into a certain type of output, then the crucial thing to work out is how the input and output are represented in the brain, and what the basic transformations are. If this can be done, then the only obstacles to building neuroprostheses are technological.

In fact, some types of neuroprostheses are already widely used. Cochlear implants can restore hearing to individuals with hearing problems – even to the profoundly deaf. They work by providing direct electrical stimulation to the auditory nerve (doing the job that would otherwise be done by hair cells in the cochlea, which is in the inner ear). This is an input prosthetic. Output prosthetics are much more complicated.

Scientists and engineers in Theodore Berger's lab at the University of Southern California have been working on a pioneering example of an output prosthetic. They have been designing and building a prosthetic implant to restore normal functioning when the hippocampus is damaged (the hippocampus plays an important role in forming and storing memories). The aim is to develop a device that will measure electrical inputs to the hippocampus; calculate what outputs would typically be generated in normal subjects, and then stimulate areas of the hippocampus to mimic a normally functioning brain. When the second edition of this afterword was being written (August 2013), an early prototype hippocampal prosthetic had been tested in rats (Berger et al. 2012) and monkeys. Now, 5 years later, the first tests have been carried out on humans, showing a significant improvement in episodic memory in epilepsy patients (Hampson et al. 2018).

Exoskeletons (robot suits) are an even more dramatic example of output prosthetics. Neuroscientists, working together with biomechanical engineers, have produced motor prostheses that restore some movement to paralyzed patients. At the opening ceremony of the 2014 Soccer World Cup in Brazil, a young Brazilian paraplegic wearing an exoskeleton got out of his wheelchair and kicked the ball to demonstrate the rehabilitation possibilities. The exoskeleton was designed by a team led by Gordon Cheng at the Technical University in Munich. The hydraulically powered suit is built from lightweight alloys. The user's EEG waves are read and translated into signals that control the robot's limbs. At the same time, feedback is sent to the user's brain from sensors in the soles of the robot's feet. The brain-computer interface making this all possible was developed by Miguel Nicolelis at the Duke University Medical Center. For an open-access review of brain-computer interfaces in rehabilitative medicine, see Lazarou et al. 2018.



## 17.4 Cognitive Science and the Law

The intensely interdisciplinary nature of cognitive science has been a recurring theme in this book. We have looked at how cognitive science has been molded by contributions from psychology, philosophy, neuroscience, linguistics, computer science, and mathematics – to give just a partial list. But the list of disciplines to which cognitive science is potentially relevant is even longer. One area where the dialog with cognitive science is gaining momentum is law.

There are many points of contact between cognitive science and the law. Eyewitness testimony is a good example. Eyewitness testimony is a fundamental pillar of almost every legal system, including retrospective identification through line-ups and similar techniques. Yet there is strong evidence that eyewitness testimony is both unreliable and manipulable (see Busey and Loftus 2007), leading to serious miscarriages of justice (subsequently discovered, for example, through DNA evidence). For those and other reasons, the National Academy of Sciences convened an expert panel to study how identification errors arise in eyewitness testimony.

The resulting report, *Identifying the Culprit: Assessing Eyewitness Identification*, released in 2014, made a number of recommendations for courts and for law enforcement agencies. For example,

- (i) line-ups should be administered by “blinded” managers who have no knowledge of which participant is suspected of the crime
- (ii) witness instructions should be standardized and designed to yield a consistent and conservative response

These and other recommendations were incorporated into Justice Department guidelines in 2017. See Albright 2017 for a review of the scientific issues emerging from the report and some of its specific recommendations.

Neuroscientific studies and findings are increasingly being used in court cases and this has given rise to the area known as neurolaw, where cognitive science intersects with forensic psychiatry, as well as legal practice. For a range of perspectives on neurolaw see the essays in Morse and Roskies 2013.

Some examples of questions actively being explored in this area:

- How can courts use neuroscientific findings to adjudicate questions about competence and capability that may arise, for example, in end of life issues?
- Are there prospects for developing brain-based techniques of lie detection?
- How should neuroscientific evidence be used as part of an insanity defense, or as evidence of mitigating circumstances?
- Does the use of neuroscientific evidence in courtrooms pose potential challenges to privacy rights?
- Can neuroscience be used to predict recidivism in offenders?

- Are there possibilities for neuroscience-based interventions that will diminish the risk of an offender reoffending (e.g., deep brain stimulation as a way of reducing sexual drive)?

Some of these issues remain in the realm of theory. But others have had serious practical ramifications. In *Miller v. Alabama*, decided in 2012, the US Supreme Court decided that mandatory life sentences without the possibility of parole are unconstitutional for juvenile offenders. Justice Elena Kagan, writing for the majority, cited an *amicus* brief from the American Psychological Association summarizing a range of evidence from neuroscience and developmental psychology on adolescent brain development and associated vulnerabilities. In *Montgomery v. Louisiana*, decided in 2016, the Court applied its earlier decision retroactively, potentially affecting more than 2,000 currently incarcerated individuals.

## 17.5 Autonomous Vehicles: Combining Deep Learning and Intuitive Knowledge

Research on self-driving cars has been one of the most visible products of the deep learning revolution in AI. Self-driving cars have been tested on semipublic roads and both large corporations and small start-ups are making ambitious predictions about when they will be widely available. In fact, 2018 was the target date proposed in 2015 both by Elon Musk of Tesla and by Google. But a series of widely publicized crashes, some fatal, have put something of a dampener on the enthusiasm.

As we saw in Chapter 12, deep learning algorithms have exponentially increased the power of machine learning, both with and without supervision. For many cognitive scientists and other observers, though, deep learning has its limitations. The great successes of deep learning have all been in relatively circumscribed domains. Chess and Go are obvious examples. But so too is image recognition. Image recognition is a purely passive activity. It is a matter of identifying patterns in a data set and then projecting those patterns onto new exemplars. Admittedly, pattern recognition is no mean achievement. Traditional AI and machine learning were only able to make limited progress for many years, whereas deep learning algorithms can outperform human experts on many complex tasks.

But still, many have thought that self-driving cars (and other forms of autonomous vehicles, such as submarines, planes, and drones) need more than sensitivity to patterns and the ability to learn from experience. They need to be able to deal with the unexpected – completely unpredictable behavior from other drivers, pedestrians, cyclists (not to mention wild animals and livestock). Will it always be possible for a self-driving car to deal with unpredictable situations by extrapolating from its training set? Boosters for current designs for self-driving cars think so, but others are not so sure.

Some cognitive scientists see a need for self-driving cars to have an analog of the intuitive knowledge that humans call common sense. Specifically, they observe that human drivers are constantly exploiting their knowledge of physical objects and how objects move and behave, as well as their knowledge of other drivers and road-users. In



other words, human driving exploits some of the key abilities that we have looked at in earlier chapters – what we called folk physics in Chapter 11 and mindreading in Chapters 13 and 14. So, a key challenge, perhaps *the* key challenge, for designers of self-driving cars is how to equip their vehicles with this kind of very general knowledge.

It may turn out that folk physics and mindreading abilities ultimately rest on very sophisticated forms of pattern recognition, so that learning them is within the reach of some version of deep learning algorithms (although perhaps ones that differ significantly from existing algorithms and networks). Alternatively, they may depend upon fundamentally different forms of learning. This second possibility is where an MIT spin-off company called iSee is placing its bets. The company is developing forms of probabilistic programming inspired by Bayesian models of the mind to equip self-driving cars with the common sense that they seem currently to lack.<sup>1</sup>

These are just some of the exciting challenges and opportunities opening up for cognitive scientists in the years ahead. I hope that readers of this book will pursue these – and, I hope, develop others.

<sup>1</sup> The founders of iSee are associated with Joshua Tenenbaum's lab at MIT. We looked at his work in Chapter 12. There is a short article on iSee in the online MIT Technology Review (published on September 20, 2017).



## GLOSSARY

**abduction (abductive reasoning):** a form of reasoning in which one derives a conclusion as the best explanation of given evidence, even though it is not entailed by the evidence that it explains.

**absolute judgment:** a judgment about the intrinsic properties of a stimulus (e.g., naming a color or identifying the pitch of a particular tone), as opposed to a relative judgment comparing two stimuli.

**access consciousness (or A-consciousness):** information available or “poised” for conscious thought and action.

**action potentials:** electrical impulses fired by neurons down their axons to other neurons.

**activation function:** a function that assigns an output signal to a neural network unit on the basis of the total input to that unit.

**algorithm:** a finite set of unambiguous rules that can be systematically applied to an object or set of objects to transform it or them in definite ways in a finite amount of time.

**anatomical connectivity:** the anatomical connections between different brain regions.

**anterograde amnesia:** the loss of memory of events after the onset of a brain injury.

**artificial neural network (connectionist network):** an abstract mathematical tool for modeling cognitive processes that uses **parallel processing** between intrinsically similar units (artificial neurons) organized in a single- or multilayer form.

**attractor:** a region in the state space of **dynamical systems** on which many different trajectories converge.

**autoencoder:** type of **unsupervised** neural network important in **deep learning**. It can learn to represent features by compressing data through a bottleneck.

**backpropagation algorithm:** a learning algorithm in multilayer neural networks in which error is spread backward through the network from the output units to the hidden units, allowing the network to modify the weights of the units in the hidden layers.

**Bayesianism:** movement in statistics that interprets probability subjectively and holds that rational thinkers will update their probabilities according to **Bayes's Rule** and that rational agents will maximize **expected utility**.

**Bayes's Rule:** gives a way of calculating the **posterior probability** of a hypothesis, conditional upon some evidence. To apply it, you need to know the **prior probability** of the hypothesis and the **likelihood** of the evidence, conditional upon the hypothesis.

**behavior-based robotics:** movement in robot design that moves beyond purely reactive **subsumption architectures** by allowing robots to represent their environment and to plan ahead.

**behaviorism:** the school of psychology holding that psychologists should only study observable phenomena and measurable behavior. Behaviorists maintain that all learning is the result of either **classical/Pavlovian** or **operant conditioning**.

**binding problem:** the problem of explaining how information processed in separate neural areas of the information-processing pathway is combined to form **representations** of objects.

**binocular rivalry:** phenomenon that occurs when different images are presented to each eye.

Subjects experience an alternation of the images.

**biorobotics:** the enterprise of designing and building models of biological organisms that reflect the basic design principles of those organisms.

**bit:** a measure of the information necessary to decide between two equally likely alternatives. For decisions between  $n$  alternatives, the number of bits =  $\log_2 n$ .

**blindsight:** a neurological disorder typically resulting from lesions in the primary visual cortex (V1, or the striate cortex). Like **unilateral spatial neglect** patients, blindsight patients report little to no awareness in one side of their visual field.

**BOLD signal:** the Blood Oxygen Level Dependent (BOLD) signal measures the contrast between oxygenated and deoxygenated hemoglobin in the brain, generally held to be an index of cognitive activity. The increase in blood oxygen can be detected by an fMRI scanner because oxygenated and deoxygenated hemoglobin have different magnetic properties.

**Boolean function:** a function that takes sets of truth values as input and produces a single truth value as output.

**Brodmann areas:** different regions of the cerebral cortex identified by the neurologist Korbinian Brodmann. The primary visual cortex, for example, is Brodmann area 17.

**cerebral cortex:** the parts of the brain, popularly called “gray matter,” that evolved most recently.

**channel capacity:** the maximum amount of data that an **information channel** can reliably transmit.

**chatbot:** a program set to respond to certain cues by making one of a small set of preprogrammed responses; these programs cannot use language to report on or navigate their environments because they do not analyze the syntactic structure or meaning of the sentences they encounter.

**cheater detection module:** hypothetical cognitive system specialized for identifying a “free rider” in a social exchange (i.e., a person who is reaping benefits without paying the associated costs).

**chunking:** Miller’s method of relabeling a sequence of information to increase the amount of data that the mind can reliably transmit, for example, relabeling sequences of digits with single numbers, i.e., 1100100 becomes “100.”

**Church-Turing thesis:** the thesis that the algorithmically calculable functions are exactly the functions that can be computed by a **Turing machine**.

**classical/Pavlovian conditioning:** the process of creating an association between a reflex response and an initially neutral stimulus by pairing the neutral stimulus (e.g., a bell) with a stimulus (e.g., food) that naturally elicits the response (e.g., salivation).

**competitive network:** an example of an artificial neural network that works by unsupervised learning.

**computation:** purely mechanical procedure for manipulating information.

**computational neuroscience:** the use of abstract mathematical models to study how the collective activities of a population of neurons could solve complex information-processing tasks.

**conditional probability:** the probability that some proposition A is true, on the assumption that some other proposition is true, e.g., the probability that it is raining (A) if we assume that it is cloudy (B). Written as  $p(A|B)$ .

**congruence priming:** a priming task in which the basic category of a prime (e.g., a tool) enhances the salience of other stimuli matching that category (e.g., other tools).

**connectionist network:** see **artificial neural network**.

**connectivity, anatomical:** physiological connections between segregated and distinct cortical regions.

**contralateral organization:** occurs when each hemisphere of the brain processes input information from the opposite side of space (e.g., when an auditory stimulus presented to the right ear is processed by the left hemisphere of the brain).

- convolutional neural network:** type of **deep learning** neural network very important in machine vision. Characterized by **sparse connectivity**, **shared weights**, and invariance under translation. Particularly well suited to image recognition.
- co-opted system:** according to **simulation theory**, a system specialized for a specific cognitive task that is then used to perform related mindreading tasks.
- corpus callosum:** the large bundle of fibers connecting the two hemispheres of the brain.
- counterfactual:** a statement about what would have happened had things been different.
- covert attention:** the possibility of directing attention at different peripheral areas while gaze is fixated on a central point.
- cross-lesion disconnection experiments:** experiments designed to trace connections between cortical areas in order to determine the pathways along which information flows. These experiments take advantage of the fact that the brain is divided into two hemispheres, with the major cortical areas being the same on each side.
- cross-talk:** the process in which separate subsystems collaborate in solving information-processing problems using each other's outputs as inputs.
- decision trees:** a branching representation of all possible paths through a problem space starting from an initial point.
- deep learning:** very impactful approach to machine learning using multilayered and hierarchically organized artificial neural networks to achieve **representation learning**.
- deep structure:** in Chomskyan linguistics, the deep structure of a sentence is its "real" syntactic structure, which serves as the basis for fixing its meaning. Two sentences with different **surface structures** can have the same deep structure (e.g., "John kissed Mary" and "Mary was kissed by John").
- dichotic listening experiments:** experiments in which subjects are presented with information in each ear in order to investigate selective attention in the auditory system.
- dishabituation paradigm:** a method for studying infant cognition that exploits the fact that infants look longer at events that they find surprising.
- distributed representation:** occurs when (as in many connectionist networks) objects or properties are represented through patterns of activation across populations of neurons, rather than through individual and discrete symbols.
- domain-specific:** term used to characterize cognitive mechanisms (modules) that carry out a very specific information-processing task with a fixed field of application.
- dorsal pathway:** the neural pathway believed to be specialized for visual information relevant to locating objects in space. This pathway runs from the **primary visual cortex** to the posterior parietal lobe.
- double dissociation:** experimental discovery that each of two cognitive functions can be performed independently of the other.
- dynamical systems hypothesis:** radical proposal to replace information-processing models in cognitive science with models based on the mathematical tools of dynamical systems theory.
- dynamical systems theory:** branch of applied mathematics using difference or differential equations to describe the evolution of physical systems over time.
- early selection model:** a cognitive model of attention in which attention operates as a filter early in the perceptual process and acts on low-level physical properties of the stimulus.
- EEG (electroencephalography):** experimental technique for studying the electrical activity of the brain.
- effective connectivity:** the causal flow of information between different brain regions.
- entropy:** a measure of how well a particular attribute classifies a set of examples. The closer the entropy is to 0, the better the attribute classifies the set.

**event-related fMRI:** neuroimaging technique for measuring the **BOLD** signal associated with rapidly changing neural events, which is possible because of the linear nature of the hemodynamic response.

**event-related potentials (ERPs)/event-related magnetic fields:** cortical signals that reflect neural network activity that can be recorded noninvasively using EEG or MEG.

**expected utility:** the expected utility of an action is the sum of the **utility** anticipated from each of its possible outcomes, each discounted by its probability. In **Bayesianism**, rational agents maximize expected utility.

**expert systems research:** a field of artificial intelligence that aims to reproduce the performance of human experts in a particular domain.

**false belief task:** an experimental paradigm first developed by psychologists Heinz Wimmer and Joseph Perner, exploring whether young children understand that someone might have mistaken beliefs about the world.

**feature engineering:** the process in machine learning of classifying a database in terms of relevant features.

**feature learning:** see **representation learning**.

**feedforward network:** a connectionist network in which activation spreads forward through the network; there is no spread of activation between units in a given layer or backward from one layer to the previous layer.

**fixed neural architectures:** the identification of determinate regions of the brain associated with particular types of modular processing.

**fMRI (functional magnetic resonance imaging):** technology for functional neuroimaging that measures levels of blood oxygen as an index of cognitive activity.

**folk physics:** an intuitive understanding of some of the basic principles governing how physical objects behave and interact.

**formal property:** a physical property of a representation that is not semantic (e.g., a formal property of the word “apple” is that it is composed of six letters of the English alphabet).

**fovea:** area in the center of the retina where visual acuity is highest, corresponding to the center of the visual field.

**frame problem:** the problem of developing expert systems in AI and building robots that can build into a system rules that will correctly identify what information and which inferences are relevant in a given situation.

**functional connectivity:** the statistical dependencies and correlations between activation in different brain areas.

**functional decomposition:** the process of explaining a cognitive capacity by breaking it down into subcapacities that can be analyzed separately. Each of these subcapacities can in turn be broken down into further nested subcapacities, until the process bottoms out in noncognitive components.

**functional neuroimaging:** a tool that allows brain activity to be studied noninvasively while subjects are actually performing experimental tasks (e.g., PET, fMRI).

**functional system:** a system that can be studied and understood primarily in terms of the role it plays and the task that it executes, irrespective of the mechanism of implementation. These systems are studied only at the computational level and are multiply realizable. (See **multiple realizability**.)

**global workspace theory of consciousness:** a leading theory of how mental states become conscious. According to this theory, attention makes low-level modular information available to conscious control (the “global workspace”), where the information is then “broadcast” to other areas of the brain.

- GOFAI:** good old-fashioned artificial intelligence – as contrasted, for example, with **artificial neural networks**, or **biorobotics**.
- graceful degradation:** the incremental deterioration of cognitive abilities that is imperceptible within small time frames.
- halting problem:** the problem first raised by David Hilbert of algorithmically determining whether a computer program will halt (i.e., deliver an output) for a given input.
- hard problem of consciousness:** the problem of explaining phenomenal consciousness by appealing to physical processes in the brain and using the traditional tools of cognitive science.
- Hebbian learning:** Donald Hebb's model of associative process according to which “neurons that fire together, wire together.”
- heuristic search hypothesis:** Newell and Simon's hypothesis that problems are solved by generating and algorithmically transforming symbol structures until a suitable solution structure is reached.
- hidden layer:** a layer of **hidden units** in an **artificial neural network**.
- hidden unit:** a unit (artificial neuron) in an **artificial neural network** whose inputs come from other units and whose outputs go to other units.
- informational encapsulation:** property of modular systems that operate with a proprietary database of information and are insulated from background knowledge and expectations.
- information channel:** a medium that transmits information from a sender to a receiver (e.g., a telephone cable or a neuron).
- integration, principle of:** fundamental idea of neuroscience stating that cognitive function involves the coordinated activity of networks of different brain areas, with different types of tasks recruiting different types of brain areas.
- intentional realism:** the thesis that propositional attitudes (e.g., beliefs and desires) can cause behavior.
- intentionality:** property in virtue of which symbols represent objects and properties in the world.
- interocular suppression:** a technique used to study consciousness, in which one eye is presented with an image of an object while the other eye is presented simultaneously with a high-contrast pattern that blocks conscious awareness of the presented image.
- joint visual attention:** occurs when infants look at objects, and take pleasure in doing so, because they see that another person is both looking at that object and noticing that the infant is also looking at the object.
- Knowledge Argument:** a thought experiment proposed by Frank Jackson and featuring a neuroscientist called Mary who is confined to a black-and-white room and has never experienced colors. Mary knows all the physical facts there are to be known, and yet, according to Jackson, there is a fact that she discovers when she leaves the room – the fact about what it is like for someone to see red.
- language of thought hypothesis:** a model of information processing developed by Jerry Fodor that holds that the basic symbol structures that carry information are sentences in an internal language of thought (sometimes called Mentalese) and that information processing works by transforming those sentences in the language of thought.
- late selection model:** a cognitive model of attention in which attention operates as a filter on representations of objects after basic perceptual processing is complete.
- Leibniz's Mill:** a thought experiment used by Gottfried Wilhelm Leibniz to draw a contrast between understanding the physical parts of the mind and understanding the distinctive nature of conscious perceptions.

**lexical access:** the processing involved in understanding single words.

**likelihood:** important concept for **Bayes's Rule**, measuring the probability of some evidence E, conditional upon a hypothesis H – i.e., how likely you think it is that you would see the evidence if the hypothesis were true. It is standardly written as the **conditional probability**  $p(E/H)$ .

**linear separability:** characteristic of **Boolean functions** that can be learned by neural networks using the **perceptron convergence learning rule**.

**local algorithm:** a learning algorithm in a connectionist network in which an individual unit weight changes directly as a function of the inputs to and outputs from that unit (e.g., the **Hebbian learning rule**).

**local field potential (LFP):** an electrophysiological signal believed to be correlated with the sum of inputs to neurons in a particular area.

**locus of selection problem:** the problem of determining whether attention is an **early selection** phenomenon or a **late selection** phenomenon.

**logical consequence:** a conclusion is the logical consequence of a set of premises just if there is no way of interpreting the premises and conclusion that makes the premises all true and the conclusion false.

**logical deducibility:** one formula is logically deducible from another just if there is a sequence of legitimate formal steps that lead from the second to the first.

**machine learning:** the production of an algorithm that will organize a complex database in terms of some target attribute by transforming symbol structures until a solution structure, or decision tree that will clarify incoming data, is reached.

**machine learning algorithm:** an algorithm for constructing a **decision tree** from a vast database of information.

**mandatory application:** a feature of **modules** where cognitive modules respond automatically to stimuli of the appropriate kind. They are not under any level of executive control.

**masked priming:** a priming task in which a stimulus is made invisible through presenting a second stimulus (the mask) in rapid succession.

**massive modularity hypothesis:** holds that all information processing is carried out by specialized modules that emerged in response to specific evolutionary problems (e.g., **cheater detection module**).

**MEG (magnetoencephalography):** brain imaging technique that measures electrical activity in the brain with magnetic fields.

**mental architecture:** a model of the mind as an information processor that answers the following three questions: In what format is information carried in a cognitive system? How is information in the cognitive system transformed? How is the mind organized to function as an information processor?

**metarepresentation:** metarepresentation occurs when a **representation** is used to represent another representation rather than to represent the world (e.g., a representation of another person's mental state).

**micro-world:** an artificially restrictive domain used in AI in which all objects, properties, and events are defined in advance.

**mirror neurons:** neurons in monkeys that fire both when the monkey performs a specific action and when it observes that action being performed by an observer.

**module:** cognitive system dedicated to performing a domain-specific information-processing task. Typically held to be informationally encapsulated but not necessarily to have a fixed neural architecture.

**morphological computation:** a research program in robotics for minimizing the amount of computational control required in a robot by building as much as possible of the computation directly into its physical structure.

- multilayer network:** an **artificial neural network** containing one or more **hidden layers**.
- multiple realizability:** a characteristic of functional systems whose tasks can be performed by a number of different physical manifestations. For example, a heart, when viewed as a functional system, is multiply realizable because human hearts and mechanical hearts can perform the same function.
- neuroeconomics:** interdisciplinary area where concepts and tools from economics are used to illuminate brain functioning.
- neurotransmitters:** neurochemicals that are transmitted across **synapses** in order to relay, amplify, and modulate signals between a neuron and another cell.
- object permanence:** the knowledge that an object exists even when it is not being perceived – an important milestone in children’s development.
- operant conditioning:** a type of conditioning in which an action (e.g., pushing a lever) is reinforced by a reward (e.g., food).
- overregularization errors:** systematic mistakes that children make during the process of language acquisition as they begin to internalize basic grammar rules. Children apply rules (such as adding the suffix “-s” to nouns to make them plural) to words that behave irregularly (e.g., saying “foots” instead of “feet”).
- paired-image subtraction paradigm:** an experimental technique that allows neuroimagers to identify the brain activation relevant to a particular task by filtering out activation associated with other tasks.
- parallel processing:** simultaneous activation of units in an **artificial neural network** that causes a spread of activation through the layers of the network.
- perceptron:** a single-unit (or single-layer) **artificial neural network**.
- perceptron convergence rule (delta rule):** a learning algorithm for perceptrons (single-unit networks). It changes a perceptron’s threshold and weights as a function of the difference between the unit’s actual and intended output.
- PET (positron emission tomography):** a functional neuroimaging technique in which localization of cognitive activity is identified by measuring blood flow to specific areas of the brain.
- phenomenal consciousness (or P-consciousness):** the experiential or “what it’s like” aspect of consciousness (e.g., the distinctive experience of smelling a rose or touching a piece of velvet cloth).
- phrase structure grammar:** describes the syntactic structure of a natural language sentence in terms of categories such as verb phrase and noun phrase. Permissible combinations of syntactic categories are given by phrase structure rules, e.g., the rule stating that every sentence must contain both a verb phrase and a noun phrase.
- physical symbol system:** a set of symbols (physical patterns) that can be combined to form complex symbol structures and contains processes for manipulating symbol structures. These processes can themselves be represented by symbols and symbol structures within the system.
- physical symbol system hypothesis:** Newell and Simon’s hypothesis that a physical symbol system has the necessary and sufficient means for general intelligent action.
- posterior probability:** this is the result of applying **Bayes’s Rule**. It is the probability of a hypothesis H, conditional upon some evidence E (i.e., the conditional probability  $p(H|E)$ ).
- poverty of stimulus argument:** maintains that certain types of knowledge must be innate, as they are too complicated to be learned from the impoverished stimuli to which humans are exposed (e.g., Chomsky’s argument for Universal Grammar).
- pragmatics:** the branch of linguistics concerned with the practical implication of language and what is actually communicated in a given context.

**predicate calculus:** formal system for exploring the logical relations between formulas built up from symbols representing individuals, properties, and logical operations. Unlike the **propositional calculus**, the predicate calculus includes quantifiers (ALL or SOME) that allow representations of generality.

**prestriate cortex:** an area in the occipital and parietal lobes that receives cortical output from the **primary visual cortex**.

**primary visual cortex:** the point of arrival in the cortex for information from the retina; also called the striate cortex and Brodmann area 17.

**priming:** an experimental technique, particularly useful in studying consciousness, where a stimulus (often not consciously perceived) influences performance on subsequent tasks.

**principle of cohesion:** principle of infant **folk physics** according to which two surfaces are part of the same object if and only if they are in contact.

**principle of contact:** principle of infant **folk physics** according to which only surfaces that are in contact can move together.

**principle of continuity:** principle of infant **folk physics** according to which objects can only move on a single continuous path through space-time.

**principle of solidity:** principle of infant **folk physics** according to which there cannot be more than one object in a place at one time.

**prior probability:** in applying **Bayes's Rule**, this is the probability assigned to the hypothesis before taking into account the evidence.

**prisoner's dilemma:** any social exchange interaction between two players where a player benefits most if she defects while her opponent cooperates and suffers most when she cooperates and her opponent defects. If each player is rational and works backward from what her opponent might do, she will always reason that the best choice is to defect.

**propositional attitude:** a psychological state that can be analyzed into a proposition (e.g., the proposition that it is snowing in St. Louis) and an attitude to that proposition (e.g., the attitude of belief or the attitude of hope).

**propositional calculus:** formal system for exploring the logical relations between formulas built up from symbols for complete propositions using logical operators (such as NOT, OR, and AND).

**psychophysics:** the branch of psychology that studies the relationship between physical stimuli and how subjects perceive and discriminate them.

**recurrent network:** an **artificial neural network** that has a feedback loop serving as a memory of what the **hidden layer** was doing at the previous time step.

**recursive definition:** process for defining a set of objects by starting with a set of base cases and specifying which transformations of objects preserve membership in the set. So, for example, a recursive definition of a **well-formed formula** in the **propositional calculus** starts with propositional symbols (the base cases) and indicates which logical operations (e.g., negation) create new formulas from existing formulas.

**reduction:** the process of showing how higher-level parts of science (e.g., thermodynamics) can be understood in terms of more basic parts of science (e.g., statistical mechanics).

**reinforcement learning:** learning where the feedback is a reward signal (as opposed to the error signal characteristic of supervised learning).

**representation:** structure carrying information about the environment. Representations can be physical symbol structures or distributed states of neural networks.

**representation learning:** the subfield of machine learning dedicated to designing algorithms that will do their own **feature engineering** on raw data.

**retrograde amnesia:** the loss of memory of events before a brain injury.

**robot reply (to the Russian room argument):** a response to John Searle's thought experiment that claims that the **Russian room** is not intelligent because it is incapable of interacting with other Russian speakers rather than because of any gap between syntax and semantics.

**Russian room argument:** John Searle's thought experiment that attempts to refute the **physical symbol system hypothesis** by showing that there can be syntactic symbol manipulation without any form of intelligence or understanding.

**saccadic eye movements:** quick and unconscious eye movements scanning the visual field.

**segregation, principle of:** fundamental principle of neuroscience stating that the cerebral cortex is divided into separate areas with distinct neuronal populations.

**selection processor:** mechanism hypothesized by Leslie enabling people to inhibit the default setting of a true belief. It is not until the selection processor is fully in place that children can pass the **false belief task**, according to Leslie.

**selective attention:** the ability of individuals to orient themselves toward, or process information from, only one stimulus within the environment, to the exclusion of others.

**semantic priming:** a priming task in which the priming effect is due to the meaning of words and not due to their phonology (how they are pronounced) or their orthography (how they are spelled).

**semantic property:** a property of a representation that holds in virtue of its content, i.e., how it represents the world (e.g., a semantic property of the word "apple" is the fact that it represents a crisp and juicy fruit).

**shared weights:** feature of **artificial neural networks** where multiple units in a single layer have the same weights. Makes for more efficient processing.

**simulation theory (radical):** the theory that mindreading takes place when people think about the world from another person's perspective rather than thinking about the other person's psychological states.

**simulation theory (standard):** the theory that people are able to reason about the mental states of others and their consequent potential behaviors by inputting "pretend" beliefs and desires into their own decision-making systems.

**situated cognition:** situated cognition theorists complain that traditional cognitive science has focused on disembodied systems that operate in highly simplified and prepackaged environments. They call instead for an approach to cognitive science that takes seriously the fact that cognitive agents are both embodied and situated within a complex environment.

**sparse connectivity:** feature of **artificial neural networks** where units in a given layer are only connected to a proper subset of units in the next layer.

**spatial resolution:** the degree of spatial detail provided by a particular technique for studying the brain.

**state space:** the state space of a system is a geometrical representation of all the possible states that the system can be in. It has as many dimensions, as the system has independently varying quantities.

**striate cortex:** see **primary visual cortex**.

**subcortex:** the part of the brain, popularly called "white matter," that developed earlier in evolution than the cerebral cortex.

**subsumption architecture:** architectures in robotics that are built up incrementally from semi-autonomous layers. Subsumption architectures (originally proposed by Rodney Brooks) typically exploit direct links between perception and action.

**supervised learning:** learning (e.g., in neural networks) that involves an explicit error signal.

**surface structure:** in Chomskyan linguistics, the surface structure of a sentence is given by the actual arrangement of written or spoken lexical items – as opposed to its **deep structure**.

- symbol-grounding problem:** the problem of determining how syntactically manipulated symbols gain semantic meaning.
- synapse:** the site where the end of an axon branch comes close to a dendrite or the cell body of another neuron. This is where signals are transmitted from one neuron to another.
- systems neuroscience:** the investigation of the function of neural systems, such as the visual system or auditory system.
- systems reply (to the Chinese room argument):** a response to John Searle's thought experiment claiming that the Chinese room as a whole understands Chinese, even though the person inside the room does not.
- temporal resolution:** the degree of temporal detail provided by a particular technique for studying the brain.
- theory of mind mechanism (TOMM):** a hypothesized cognitive system specialized for attributing **propositional attitudes** and using those attributions to predict and explain behavior.
- threshold:** the minimum amount of activity necessary to initiate the firing of a unit in an **artificial neural network**.
- TIT FOR TAT:** a successful strategy used in social exchanges, such as the **prisoner's dilemma** whereby a player cooperates with his opponent during the first round and in subsequent rounds copies the action taken by the opponent on the preceding round.
- transformational grammar:** a theoretical account of the rules governing how surface structures in natural languages are generated from **deep structures**.
- truth condition:** the state of affairs that makes a particular statement true.
- truth rule:** a rule that states the truth condition for a given statement.
- Turing machine:** a theoretical model of an abstract computation device that can (according to the Church-Turing thesis) compute any effectively calculable function.
- unilateral spatial neglect:** a neurological disorder typically due to damage to the posterior parietal cortex in one hemisphere in which patients describe themselves as unaware of stimuli in the contralateral half of their visual field.
- unsupervised learning:** learning (e.g., in neural networks) where there is no explicit error or reward signal.
- utility:** a widely used concept for measuring the value of an action or outcome to an individual. In **Bayesianism**, rational agents maximize **expected utility**.
- ventral pathway:** the neural pathway believed to be specialized for visual information relevant to recognizing and identifying objects. This pathway runs from the primary visual cortex to the temporal lobe.
- Wason selection task:** experiment developed to test people's understanding of conditional reasoning. Subjects are asked to identify the additional information they would need in order to tell if a given conditional statement is true or false.
- well-formed formula:** a string of symbols in a formal language that is legitimately constructed through the formation rules of that language.



## BIBLIOGRAPHY

- Abu-Akel, A., and Shamay-Tsoory, S. (2011). Neuroanatomical and neurochemical bases of theory of mind. *Neuropsychologia*, 49, 2971–84.
- Adams, F., and Aizawa, A. (2010). *The Bounds of Cognition*. Oxford: Wiley-Blackwell.
- Adolphs, R. (2009). The social brain: Neural basis of social knowledge. *Annual Review of Psychology*, 60, 693–716.
- Adolphs, R., and Tranel, D. (2000). Emotion recognition and the human amygdala. In J. P. Aggleton (ed.), *The Amygdala: A Functional Analysis*. Oxford: Oxford University Press.
- Adolphs, R., Tranel, D., Damasio, H., and Damasio, A. (1994). Impaired recognition of emotion in facial expressions following bilateral damage to the human amygdala. *Nature*, 372, 669–72.
- Aglioti, S., DeSouza, J. F. X., and Goodale, M. A. (1995). Size-contrast illusions deceive the eye but not the hand. *Current Biology*, 5, 679–85.
- Albright, T. D. (2017). Why eyewitnesses fail. *Proceedings of the National Academy of Sciences*, 114 (30), 7758–64.
- Anderson, J. A. (2003). McCulloch-Pitts neurons. In L. Nadel (ed.), *Encyclopedia of Cognitive Science*. New York: Nature Publishing Group.
- Anderson, J. R., Bothell, D., Byrne, M. D., Douglass, S., Lebiere, C., and Qin, Y. (2004). An integrated theory of the mind. *Psychological Review*, 4, 1036–160.
- Anderson, M. L., Richardson, M. J., and Chemero, A. (2012). Eroding the boundaries of cognition: Implications of embodiment. *Topics in Cognitive Science*, 4, 717–30.
- Apperly, I. A., Samson, D., and Humphreys, G. W. (2005). Domain-specificity and theory of mind: Evaluating neuropsychological evidence. *Trends in Cognitive Sciences*, 9, 572–7.
- Apperly, I. A., Samson, D., Chiavarino, C., and Humphreys, G. W. (2004). Frontal and temporoparietal lobe contributions to theory of mind: Neuropsychological evidence from a false-belief task with reduced language and executive demands. *Journal of Cognitive Neuroscience*, 16, 1773–84.
- Arbib, M. A. (1987). *Brains, Machines, and Mathematics*. New York: Springer.
- Arbib, M. A. (2003). *The Handbook of Brain Theory and Neural Networks*. Cambridge, MA: MIT Press.
- Arkin, R. C. (1998). *Behavior-Based Robotics*. Cambridge, MA: MIT Press.
- Ashby, F. G. (2011). *Statistical Analysis of fMRI Data*. Cambridge, MA: MIT Press.
- Aslin, R. N., Saffran, J. R., and Newport, E. L. (1998). Computation of conditional probability statistics by 8-month-old infants. *Psychological Science*, 9, 321–4.
- Baars, B. J. (1988). *A Cognitive Theory of Consciousness*. Cambridge: Cambridge University Press.
- Baars, B. J. (2002). The conscious access hypothesis: Origins and recent evidence. *Trends in Cognitive Sciences*, 6, 47–52.
- Baars, B. J., and Gage, N. M. (eds.) (2010). *Cognition, Brain, and Consciousness: An Introduction to Cognitive Neuroscience* (2nd edn.). Burlington, MA: Elsevier.
- Baars, B. J., and Gage, N. M. (2012). *Fundamentals of Cognitive Neuroscience: A Beginner's Guide*. Waltham, MA: Academic Press.
- Baddeley, A. D. (2003). Working memory: Looking back and looking forward. *Nature Reviews Neuroscience*, 4, 829–39.

- Baddeley, A. D. (2007). *Working Memory, Thought, and Action*. New York: Oxford University Press.
- Baddeley, A. D., and Hitch, G. J. L. (1974). Working memory. In G. A. Bower (ed.), *The Psychology of Learning and Motivation: Advances and Research*. New York: Academic Press.
- Baillargeon, R. (1986). Representing the existence and the location of hidden objects: Object permanence in 6- and 8-month-old infants. *Cognition*, 23, 21–41.
- Baillargeon, R. (1987). Object permanence in 3- and 4-month-old infants. *Developmental Psychology*, 23, 655–64.
- Baillargeon, R., and Carey, S. (2012). Core cognition and beyond: The acquisition of physical and numerical knowledge. In S. Pauen (ed.), *Early Childhood Development and Later Outcome*. Cambridge: Cambridge University Press.
- Baillargeon, R., Scott, R. M., and He, Z. (2010). False-belief understanding in infants. *Trends in Cognitive Sciences*, 14, 110–18.
- Baillargeon, R., Li, J., Gertner, Y., and Wu, D. (2010). How do infants reason about physical events? In U. Goswami (ed.), *The Wiley-Blackwell Handbook of Childhood Cognitive Development* (2nd edn.). Oxford: Blackwell.
- Bandettini, P. A., and Ungerleider, L. G. (2001). From neuron to BOLD: New connections. *Nature Neuroscience*, 4, 864–6.
- Baron-Cohen, S. (1995). *Mindblindness: An Essay on Autism and Theory of Mind*. Cambridge, MA: MIT Press.
- Baron-Cohen, S. (2005). The empathizing system: A revision of the 1994 model of the mindreading system. In B. Ellis and D. Bjorklund (eds.), *Origins of the Social Mind*. New York: Guilford.
- Baron-Cohen, S. (2009). The empathizing-systemizing (E-S) theory. *Annals of the New York Academy of Sciences*, 1156, 68–80.
- Baron-Cohen, S., Leslie, A. M., and Frith, U. (1985). Does the autistic child have a “theory of mind”? *Cognition*, 21, 37–46.
- Baron-Cohen, S., Tager-Flusberg, H., and Cohen, D. J. (eds.) (2000). *Understanding Other Minds: Perspectives from Developmental Cognitive Neuroscience*. New York: Oxford University Press.
- Barrett, H. C., and Kurzban, R. (2006). Modularity in cognition: Framing the debate. *Psychological Review*, 113, 628–47.
- Bassett, D. S., and Bullmore, E. (2006). Small-world brain networks. *Neuroscientist*, 12, 512–23.
- Bayne, T. (2012). *The Unity of Consciousness*. New York: Oxford University Press.
- Beate, S. (2011). Theory of mind in infancy. *Child Development Perspectives*, 5, 39–43.
- Bechtel, W. (1999). Unity of science. In R. A. Wilson and F. Keil (eds.), *The MIT Encyclopedia of Cognitive Science*. Cambridge, MA: MIT Press.
- Bechtel, W., and Abrahamsen, A. A. (2002). *Connectionism and the Mind: Parallel Processing, Dynamics and Evolution in Networks*. Cambridge, MA: Blackwell.
- Bechtel, W., Mandik, P., Mundale, J., and Stufflebeam, R. S. (eds.) (2001). *Philosophy and the Neurosciences: A Reader*. Malden, MA: Blackwell.
- Berger, T. W., Song, D., Chan, R. H. M., Marmarelis, V. Z., LaCoss, J., Wills, J., ... Granacki, J. J. (2012). A hippocampal cognitive prosthesis: Multi-input, multi-output nonlinear modeling and VLSI implementation. *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, 20(2), 198–211.
- Bermúdez, J. L. (2005). *Philosophy of Psychology: A Contemporary Introduction*. New York: Routledge.
- Bermúdez, J. L. (ed.) (2006). *Philosophy of Psychology: Contemporary Readings*. London: Routledge.
- Bermúdez, J. L. (2009). *Decision Theory and Rationality*. Oxford: Oxford University Press.
- Bernstein, I. H., Bissonnette, V., Vyas, A., and Barclay, P. (1989). Semantic priming: Subliminal perception or context? *Perception and Psychophysics*, 45, 153–61.

- Berti, A., and Rizzolatti, G. (1992). Visual processing without awareness: Evidence from unilateral neglect. *Journal of Cognitive Neuroscience*, 4, 345–51.
- Berwick, R. C., Pietroski, P., Yankama, B., and Chomsky, N. (2011). Poverty of the stimulus revisited. *Cognitive Science*, 35(7), 1207–42.
- Bickle, J. (2006). Reducing mind to molecular pathways: Explicating the reductionism implicit in current cellular and molecular neuroscience. *Synthese*, 151, 411–34.
- Bisiach, E., and Luzzatti, C. (1978). Unilateral neglect of representational space. *Cortex*, 14, 129–33.
- Blamire, A. M., Ogawa, S., Ugurbil, K., Rothman, D., McCarthy, G., Ellermann, J. M., ... Shulman, R. G. (1992). Dynamic mapping of the human visual cortex by high-speed magnetic resonance imaging. *Proceedings of the National Academy of Sciences*, 89(22), 11069–73.
- Blank, A. (2010). On interpreting Leibniz's Mill. In P. K. Machamer and G. Wolters (eds.), *Interpretation: Ways of Thinking about the Sciences and the Arts*. Pittsburgh, PA: University of Pittsburgh Press.
- Block, N. (ed.) (1981). *Imagery*. Cambridge, MA: MIT Press.
- Block, N. (1995a). The mind as the software of the brain. In D. Osherson, L. Gleitman, S. M. Kosslyn, E. Smith, and R. J. Sternberg (eds.), *An Invitation to Cognitive Science*. Cambridge, MA: MIT Press.
- Block, N. (1995b). On a confusion about the function of consciousness. *Behavioral and Brain Sciences*, 18, 227–47.
- Block, N. (2007). Consciousness, accessibility, and the mesh between psychology and neuroscience. *Behavioral and Brain Sciences*, 30, 481–548.
- Block, N. (2011). Perceptual consciousness overflows cognitive access. *Trends in Cognitive Sciences*, 15, 567–75.
- Block, N., Flanagan, O., and Güzeldere, G. (eds.) (1997). *The Nature of Consciousness: Philosophical Debates*. Cambridge, MA: MIT Press.
- Bloom, P., and German, T. P. (2000). Two reasons to abandon the false belief task as a test of theory of mind. *Cognition*, 77, B25–B31.
- Boden, M. A. (1977). *Artificial Intelligence and Natural Man*. Brighton: Harvester Press.
- Boden, M. A. (1990a). Escaping from the Chinese room. In *The Philosophy of Artificial Intelligence*. Oxford: Oxford University Press.
- Boden, M. A. (ed.) (1990b). *The Philosophy of Artificial Intelligence*. Oxford: Oxford University Press.
- Boden, M. A. (2006). *Mind as Machine: A History of Cognitive Science*. Oxford: Oxford University Press.
- Bornhofen, C., and McDonald, S. (2008). Emotion perception deficits following traumatic brain injury: A review of the evidence and rationale for intervention. *Journal of the International Neuropsychological Society*, 14(4), 511–25.
- Boucher, J. (1996). What could possibly cause autism? In P. Carruthers and P. K. Smith (eds.), *Theories of Theory of Mind*. Cambridge: Cambridge University Press.
- Bowers, J. S. (2009). On the biological plausibility of grandmother cells: Implications for neural network theories in psychology and neuroscience. *Psychological Review*, 16, 220–51.
- Brachman, R. J., and Levesque, H. J. (eds.) (1985). *Readings in Knowledge Representation*. Los Altos, CA: M. Kaufmann.
- Bremner, G. J. (1994). *Infancy*. Oxford: Wiley-Blackwell.
- Bressler, S. L., Tang, W., Sylvester, C. M., Shulman, G. L., and Corbetta, M. (2008). Top-down control of human visual cortex by frontal and parietal cortex in anticipatory visual spatial attention. *Journal of Neuroscience*, 28, 10056–61.
- Brewer, J. B., Zhao, Z., Desmond, J. E., Glover, G. H., and Gabrieli, J. D. E. (1998). Making memories: Brain activity that predicts how well visual experience will be remembered. *Science*, 281, 1185–7.

- Broadbent, D. E. (1954). The role of auditory localization in attention and memory span. *Journal of Experimental Psychology*, 47, 191–6.
- Broadbent, D. E. (1958). *Perception and Communication*. London: Pergamon Press.
- Brook, A. (2007). *The Prehistory of Cognitive Science*. New York: Palgrave Macmillan.
- Brooks, R. (1991). Intelligence without representation. *Artificial Intelligence*, 47, 139–59. Reprinted in J. Haugeland (ed.) (1997), *Mind Design II: Philosophy, Psychology, Artificial Intelligence*. Cambridge, MA: MIT Press.
- Brooks, R. (1999). *Cambrian Intelligence: The Early History of the New AI*. Cambridge, MA: MIT Press.
- Broome, J. (1991). Utility. *Economics and Philosophy* 7: 1–12.
- Buckley, R. F., Schultz, A. P., Hedden, T., Papp, K. V., Hanseeuw, B. J., Marshall, G., ... Chhatwal, J. P. (2017). Functional network integrity presages cognitive decline in preclinical Alzheimer disease. *Neurology*, 89(1), 29–37.
- Buckner, R. L. (1998). Event-related fMRI and the hemodynamic response. *Human Brain Mapping*, 6(5–6), 373–7.
- Bullmore, E., and Sporns, O. (2009). Complex brain networks: Graph theoretical analysis of structural and functional systems. *Nature Reviews Neuroscience*, 10, 186–98.
- Busey, T. A., and Loftus, G. R. (2007). Cognitive science and the law. *Trends in Cognitive Sciences*, 11(3), 111–17.
- Byrne, R. M. J., and Johnson-Laird, P. N. (2009). "If" and the problems of conditional reasoning. *Trends in Cognitive Sciences*, 13, 282–7.
- Cacchione, T. (2013). The foundations of object permanence: does perceived cohesion determine infants' appreciation of the continuous existence of material objects? *Cognition*, 128(3), 397–406.
- Cadieu, C. F., Hong, H., Yamins, D. L. K., Pinto, N., Ardila, D., Solomon, E. A., ... DiCarlo, J. J. (2014). Deep neural networks rival the representation of primate IT cortex for core visual object recognition. *PLoS Computational Biology*, 10(12), e1003963.
- Calvo, P., and Symons, J. (2014). *The Architecture of Cognition: Rethinking Fodor and Pylyshyn's Systematicity Challenge*. Cambridge, MA: MIT Press.
- Carey S. (2009). *The Origin of Concepts*. Oxford: Oxford University Press.
- Carey, S., and Spelke, E. S. (1996). Science and core knowledge. *Philosophy of Science*, 63, 515–33.
- Carrasco, M. (2011). Visual attention: The past 25 years. *Vision Research*, 51, 1484–535.
- Carrasco, M. (2018). How visual spatial attention alters perception. *Cognitive Processing*, 19(Suppl. 1), 77–88.
- Carrington, S. J., and Bailey, A. J. (2009). Are there theory of mind regions in the brain? A review of the neuroimaging literature. *Human Brain Mapping*, 30, 2313–35.
- Carruthers, P. (2000). *Phenomenal Consciousness*. Cambridge: Cambridge University Press.
- Carruthers, P. (2006). *The Architecture of the Mind*. Cambridge: Cambridge University Press.
- Carruthers, P. (2008a). On Fodor-fixation, flexibility, and human uniqueness: A reply to Cowie, Machery, and Wilson. *Mind and Language*, 23, 293–303.
- Carruthers, P. (2008b). Precis of *The Architecture of the Mind: Massive Modularity and the Flexibility of Thought*. *Mind and Language*, 23, 257–62.
- Carruthers, P. (2013). Mindreading in infancy. *Mind and Language*, 28, 141–72.
- Carruthers, P., and Smith, P. K. (eds.) (1996). *Theories of Theory of Mind*. Cambridge: Cambridge University Press.
- Chalmers, D. (1995). Facing up to the problem of consciousness. *Journal of Consciousness Studies*, 2, 200–19.
- Chalmers, D. (1996). *The Conscious Mind*. Oxford: Oxford University Press.
- Chalmers, D. J. (2013). How can we construct a science of consciousness? *Annals of the New York Academy of Sciences*, 1303, 25–35.

- Charpac, S., and Stefanovic, B. (2012). Shedding light on the BOLD fMRI response. *Nature Methods*, 9, 547–9.
- Chater, N., and Manning, C. D. (2006). Probabilistic models of language processing and acquisition. *Trends in Cognitive Sciences*, 10(7), 335–44.
- Chater, N., and Oaksford, M. (2008). *The Probabilistic Mind: Prospects for Bayesian Cognitive Science*. Oxford: Oxford University Press.
- Chater, N., Oaksford, M., Hahn, U., and Heit, E. (2010). Bayesian models of cognition. *Wiley Interdisciplinary Reviews: Cognitive Science*, 1(6): 811–23.
- Chelazzi, L., and Corbetta, M. (2000). Cortical mechanisms of visuospatial attention in the primate brain. In M. S. Gazzaniga (ed.), *The New Cognitive Neurosciences* (2nd edn.). Cambridge, MA: MIT Press.
- Chemero, A. (2009). *Radical Embodied Cognitive Science*. Cambridge, MA: MIT Press.
- Cherry, E. C. (1953). Some experiments on the recognition of speech, with one and two ears. *Journal of the Acoustical Society of America*, 25, 975–9.
- Chhatwal, J. P., Schultz, A. P., Johnson, K., Benzinger, T. L. S., Jack, C., Ances, B. M., . . . Sperling, R. A. (2013). Impaired default network functional connectivity in autosomal dominant Alzheimer disease. *Neurology*, 81(8), 736–44.
- Chomsky, N. (1957). *Syntactic Structures*. Gravenhage: Mouton.
- Chomsky, N. (1959). A review of B. F. Skinner's *Verbal Behavior*. *Language*, 35, 26–58.
- Chomsky, N. (1968). *Language and Mind*. New York: Harper and Row.
- Chomsky, N. (1980a). *Rules and Representations*. New York: Columbia University Press.
- Chomsky, N. A. (1980b). Rules and representations. *Behavioral and Brain Sciences*, 3(127), 1–61.
- Christiansen, M. H., and Chater, N. (2001). *Connectionist Psycholinguistics*. Westport, CT: Ablex.
- Chun, M. M., Golomb, J. D., and Turk-Browne, N. B. (2011). A taxonomy of external and internal attention. *Annual Review of Psychology*, 62, 73–101.
- Churchland, P. M. (1990a). On the nature of theories: A neurocomputational perspective. In *A Neurocomputational Perspective: The Nature of Mind and the Structure of Science*. Cambridge, MA: MIT Press.
- Churchland, P. M. (1990b). Cognitive activity in artificial neural networks. In N. Block and D. Osherson (eds.), *Invitation to Cognitive Science*. Cambridge, MA: MIT Press. Reprinted in R. Cummins and D. D. Cummins (2000), *Minds, Brains, and Computers: The Foundations of Cognitive Science: An Anthology*. Malden, MA: Blackwell.
- Churchland, P. M. (2007). *Neurophilosophy at Work*. Cambridge: Cambridge University Press.
- Churchland, P. M., and Churchland, P. S. (1990). Could a machine think? *Scientific American*, 262(1), 32–7.
- Churchland, P. S. (1986). *Neurophilosophy: Toward a Unified Science of the Mind/Brain*. Cambridge, MA: MIT Press.
- Churchland, P. S., and Sejnowski, T. J. (1992). *The Computational Brain*. Cambridge, MA: MIT Press.
- Clancey, W. J. (1997). *Situated Cognition: On Human Knowledge and Computer Representations*. Cambridge: Cambridge University Press.
- Clark, A. (1989). *Microcognition: Philosophy, Cognitive Science, and Parallel Distributed Processing*. Cambridge, MA: MIT Press.
- Clark, A. (1993). *Associative Engines: Connectionism, Concepts, and Representational Change*. Cambridge, MA: MIT Press.
- Clark, A. (1997). *Being There: Putting Brain, Body, and World Together Again*. Cambridge, MA: MIT Press.
- Clark, A. (1998). Time and mind. *Journal of Philosophy*, 95, 354–76.

- Clark, A. (2001). *Mindware: An Introduction to the Philosophy of Cognitive Science*. New York: Oxford University Press.
- Clark, A. (2008). *Supersizing the Mind: Embodiment, Action, and Cognitive Extension*. New York: Oxford University Press.
- Clark, A. (2011). Précis of *Supersizing the Mind: Embodiment, Action, and Cognitive Extension*. *Philosophical Studies*, 152, 413–16.
- Clark, A. (2016). *Surfing Uncertainty: Prediction, Action, and the Enbodyed Mind*. Oxford: Oxford University Press.
- Clearfield, M. W., Dineva, E., Smith, L. B., Diedrich, F. J., and Thelen, E. (2009). Cue salience and infant perseverative reaching: Tests of the dynamic field theory. *Developmental Science*, 12, 26–40.
- Colby, C. L., and Goldberg, M. E. (1999). Space and attention in parietal cortex. *Annual Review of Neuroscience*, 22, 319–49.
- Cook, R., Bird, G., Catmur, C., Press, C., and Heyes, C. (2014). Mirror neurons: From origin to function. *Behavioral and Brain Sciences*, 37(2), 177–92.
- Cook, V. J., and Newson, M. (2007). *Chomsky's Universal Grammar: An Introduction* (3rd edn.). Oxford: Blackwell.
- Cooper, L. A., and Shepard, R. N. (1973). The time required to prepare for a rotated stimulus. *Memory and Cognition*, 1, 246–50.
- Copeland, J. G. (1993). *Artificial Intelligence: A Philosophical Introduction*. Oxford: Blackwell.
- Corbetta, M., and Shulman, G. L. (2011). Spatial neglect and attention networks. *Annual Review of Neuroscience*, 34, 569–99.
- Corkin, S. (2002). What's new with the amnesic patient H. M.? *Nature Reviews Neuroscience*, 3, 153–60.
- Cosmides, L. (1989). The logic of social exchange: Has natural selection shaped how humans reason? Studies with the Wason selection task. *Cognition*, 31, 187–276.
- Cosmides, L., and Tooby, J. (1992). Cognitive adaptations for social exchange. In J. Berkow, L. Cosmides, and J. Tooby (eds.), *The Adapted Mind: Evolutionary Psychology and the Generation of Culture*. New York: Oxford University Press.
- Cosmides, L., and Tooby, J. (1994). Origins of domain-specificity: The evolution of functional organization. In L. A. Hirschfeld and S. F. Gelman (eds.), *Mapping the Mind: Domain Specificity in Cognition and Culture*. Cambridge: Cambridge University Press. Reprinted in J. L. Bermúdez (ed.) (2006), *Philosophy of Psychology: Contemporary Readings*. London: Routledge.
- Cosmides, L., and Tooby, J. (2013). Evolutionary psychology: New perspectives on cognition and motivation. *Annual Review of Psychology*, 64, 201–29.
- Cosmides, L., Barrett, H. C., and Tooby, J. (2010). Adaptive specializations, social exchange, and the evolution of human intelligence. *Proceedings of the National Academy of Sciences USA*, 107, 9007–14.
- Cowie, F. (1999). *What's Within? Nativism Reconsidered*. Oxford: Oxford University Press.
- Cowie, F. (2008). Us, them and it: Modules, genes, environments and evolution. *Mind and Language*, 23, 284–92.
- Crane, T. (2003). *The Mechanical Mind: A Philosophical Introduction to Minds, Machines, and Mental Representation*. London: Routledge.
- Craver, C. (2007). *Explaining the Brain: Mechanisms and the Mosaic Unity of Neuroscience*. New York: Oxford University Press.
- Crick, F., and Koch, C. (2003) A framework for consciousness. *Nature Neuroscience*, 6, 119–26.
- Croker, V., and McDonald, S. (2005). Recognition of emotion from facial expression following traumatic brain injury. *Brain Injury*, 19(10), 787–99.
- Cummins, R. (2000). "How does it work?" versus "What are the laws?" In F. C. Keil and R. A. Wilson (eds.), *Explanation and Cognition*. Cambridge, MA: MIT Press.

- Cummins, R., and Cummins, D. D. (2000). *Minds, Brains, and Computers: The Foundations of Cognitive Science: An Anthology*. Malden, MA: Blackwell.
- Cutland, N. J. (1980). *Computability: An Introduction to Recursive Function Theory*. Cambridge: Cambridge University Press.
- Dale, A. M., and Buckner, R. L. (1997). Selective averaging of rapidly presented individual trials using fMRI. *Human Brain Mapping*, 5(5), 329–40.
- Davies, M., and Stone, T. (eds.) (1995a). *Folk Psychology*. Oxford: Blackwell.
- Davies, M., and Stone, T. (eds.) (1995b). *Mental Simulation*. Oxford: Blackwell.
- Davis, M. (2000). *The Universal Computer: The Road from Leibniz to Turing*. New York: Norton.
- Davis, M. (2001). *Engines of Logic: Mathematicians and the Origin of the Computer*. New York: Norton.
- Davis, S. (1993). *Connectionism: Theory and Practice*. New York: Oxford University Press.
- Dawkins, R. (1979). Twelve misunderstandings of kin selection. *Zeitschrift für Tierpsychologie*, 51, 184–200.
- Dawson, M. R. W. (1998). *Understanding Cognitive Science*. Oxford: Blackwell.
- Dawson, M. R. W. (2004). *Minds and Machines: Connectionism and Psychological Modeling*. Oxford: Blackwell.
- Dawson, M. R. W. (2005). *Connectionism: A Hands-On Approach*. Oxford: Blackwell.
- Dayan, P., and Abbott, L. F. (2005). *Theoretical Neuroscience: Computational and Mathematical Modeling of Neural Systems*. Cambridge, MA: MIT Press.
- Dehaene, S., and Changeux, J. (2011). Experimental and theoretical approaches to conscious processing. *Neuron*, 70, 200–227.
- Dehaene, S., and Naccache, L. (2001). Towards a cognitive neuroscience of consciousness: Basic evidence and a workspace framework. *Cognition*, 79, 1–37.
- Dehaene, S., Kerszberg, M., and Changeux, J. P. (1998). A neuronal model of a global workspace in effortful cognitive tasks. *Proceedings of the National Academy of Sciences USA*, 95, 14529–34.
- Dehaene, S., Changeux, J., Naccache, L., Sackur, J., and Sergent, C. (2006). Conscious, preconscious, and subliminal processing: A testable taxonomy. *Trends in Cognitive Sciences*, 10, 204–11.
- Dehaene, S., Charles, L., King, J.-R., and Marti, S. (2014). Toward a computational theory of conscious processing. *Current Opinion in Neurobiology*, 25, 76–84.
- Dennett, D. (1984). Cognitive wheels: The frame problem in artificial intelligence. In C. Hookway (ed.), *Minds, Machines, and Evolution*. Cambridge: Cambridge University Press.
- Dennett, D. C. (1969). *Content and Consciousness*. London: Routledge & Kegan Paul.
- Dennett, D. C. (1991). *Consciousness Explained*. Boston: Little, Brown.
- Deshpande, G., and Hu, X. (2012). Investigating effective brain connectivity from fMRI data: Past findings and current issues with reference to Granger causality analysis. *Brain Connectivity*, 2(5), 235–45.
- Doherty, M. J. (1999). Selecting the wrong processor: A critique of Leslie's theory of mind mechanism-selection processor theory. *Developmental Science*, 2(1), 81–85.
- Dreyfus, H. L. (1977). *Artificial Intelligence and Natural Man*. New York: Basic Books.
- Driver, J. (2001). A selective review of selective attention research from the past century. *British Journal of Psychology*, 92(1), 53–78.
- Driver, J., and Mattingly, J. B. (1998). Parietal neglect and visual awareness. *Nature Neuroscience*, 1, 17–22.
- Driver, J., and Vuilleumier, P. (2001). Perceptual awareness and its loss in unilateral neglect and extinction. *Cognition*, 79, 39–88.
- Duncan, S. (2011). Leibniz's Mill arguments against materialism. *Philosophical Quarterly*, 62, 250–72.

- Ekstrom, A.** (2010). How and when the fMRI BOLD signal relates to underlying neural activity: The danger in dissociation. *Brain Research Review*, 62(2), 233–44.
- Eliasmith, C.** (1996). The third contender: A critical examination of the dynamicist theory of cognition. *Philosophical Psychology*, 9, 441–63.
- Elliott, M. H.** (1928). The effect of change or reward on the maze performance of rats. *University of California Publications in Psychology*, 4, 19–30.
- Elman, J. L.** (2005). Connectionist models of cognitive development: Where next? *Trends in Cognitive Sciences*, 9, 111–17.
- Elman, J. L., Bates, E. A., Johnson, M. H., and Karmiloff-Smith, A.** (1996). *Rethinking Innateness: A Connectionist Perspective on Development*. Cambridge, MA: MIT Press.
- Ernst, M. O., and Banks, M. S.** (2002). Humans integrate visual and haptic information in a statistically optimal fashion. *Nature*, 415(6870), 429–33.
- Evans, J. S. B. T., and Over, D.** (2004). *If*. Oxford: Oxford University Press.
- Fang, F., and He, S.** (2005) Cortical responses to invisible objects in the human dorsal and ventral pathways. *Nature Neuroscience*, 10, 1380–5.
- Felleman, D. J., and Van Essen, D. C.** (1991). Distributed hierarchical processing in the primate cerebral cortex. *Cerebral Cortex*, 1, 1–47.
- Finkbeiner, M., and Forster, K. I.** (2008). Attention, intention and domain-specific processing. *Trends in Cognitive Sciences*, 12, 59–64.
- Flanagan, O. J.** (1991). *The Science of the Mind*. Cambridge, MA: MIT Press.
- Fletcher-Watson, S., McConnell, F., Manola, E., and McConachie, H.** (2014). Interventions based on the Theory of Mind cognitive model for autism spectrum disorder (ASD). *Cochrane Database Systematic Reviews*, 3.
- Fodor, J.** (1975). *The Language of Thought*. Cambridge, MA: Harvard University Press.
- Fodor, J.** (1983). *The Modularity of Mind*. Cambridge, MA: MIT Press.
- Fodor, J.** (1985). Precis of *The Modularity of Mind*. *Behavioral and Brain Sciences*, 1, 1–5.
- Fodor, J.** (1987). *Psychosemantics*. Cambridge, MA: MIT Press.
- Fodor, J.** (2000). *The Mind Doesn't Work That Way: The Scope and Limits of Computational Psychology*. Cambridge, MA: MIT Press.
- Fodor, J.** (2008). *LOT 2: The Language of Thought Revisited*. Oxford: Oxford University Press.
- Fodor, J., and Pylyshyn, Z.** (1988). Connectionism and cognitive architecture: A critical analysis. *Cognition*, 28, 3–71.
- Frankish, K., and Ramsey, W. (eds.)** (2012). *The Cambridge Handbook of Cognitive Science*. Cambridge: Cambridge University Press.
- Franklin, S.** (1995). *Artificial Minds*. Cambridge, MA: MIT Press.
- Franz, V. H., Gegenfurtner, K. R., Bülthoff, H. H., and Fahle, M.** (2000). Grasping visual illusions: No evidence for a dissociation between perception and action. *Psychological Science*, 11, 20–5.
- Freyberg, J., Robertson, C. E., and Baron-Cohen, S.** (2015). Reduced perceptual exclusivity during object and grating rivalry in autism. *Journal of Vision*, 15(13), 11. doi:10.1167/15.13.11.
- Friedenberg, J., and Silverman, G.** (2006). *Cognitive Science: An Introduction to the Study of Mind*. Thousand Oaks, CA: Sage.
- Friedman, O., and Leslie, A. M.** (2007). The conceptual underpinnings of pretense: Pretending is not “behaving-as-if.” *Cognition*, 105(1), 103–24.
- Friedman, O., Neary, K. R., Burnstein, C. L., and Leslie, A. M.** (2010). Is young children’s recognition of pretense metarepresentational or merely behavioral? Evidence from 2- and 3-year-olds’ understanding of pretend sounds and speech. *Cognition*, 115(2), 314–19.
- Friston, K. J.** (2011). Functional and effective connectivity: A review. *Brain Connectivity*, 1(1), 13–36.

- Friston, K., Moran, R., and Seth, A. K. (2013). Analysing connectivity with Granger causality and dynamic causal modelling. *Current Opinion in Neurobiology*, 23(2), 172–8.
- Frith, C., and Frith, U. (2012). Mechanisms of social cognition. *Annual Review of Psychology*, 63, 287–313.
- Funt, B. V. (1980). Problem-solving with diagrammatic representations. *Artificial Intelligence*, 13, 201–30. Reprinted in R. J. Brachman and H. J. Levesque (eds.) (1985), *Readings in Knowledge Representation*. Los Altos, CA: M. Kaufmann.
- Gallistel, C. R. (1990). *The Organization of Learning*. Cambridge, MA: MIT Press.
- Gardner, H. (1985). *The Mind's New Science: A History of the Cognitive Revolution*. New York: Basic Books.
- Gazzaniga, M. S. (ed.) (1995). *The New Cognitive Neurosciences* (1st edn.). Cambridge, MA: MIT Press.
- Gazzaniga, M. S. (ed.) (2000). *The New Cognitive Neurosciences* (2nd edn.). Cambridge, MA: MIT Press.
- Gazzaniga, M. S. (ed.) (2004). *The New Cognitive Neurosciences* (3rd edn.). Cambridge, MA: MIT Press.
- Gazzaniga, M. S. (ed.) (2009). *The New Cognitive Neurosciences* (4th edn.). Cambridge, MA: MIT Press.
- Gazzaniga, M. S., and Mangun, G. (eds.) (2014). *The New Cognitive Neurosciences* (5th edn.). Cambridge, MA: MIT Press.
- Gazzaniga, M. S., Halpern, T., and Heatherton, D. (2011). *Psychological Science* (4th edn.). New York: Norton.
- Gazzaniga, M. S., Ivry, R. B., and Mangun, G. R. (2013). *Cognitive Neuroscience: The Biology of the Mind* (4th edn.). New York: Norton.
- Glass, A. L. (2016). *Cognition: A Neuroscience Approach*. Cambridge: Cambridge University Press.
- Glasser, M. F., Smith, S. M., Marcus, D. S., Andersson, J. L. R., Auerbach, E. J., Behrens, T. E. J., ... Van Essen, D. C. (2016). The Human Connectome Project's neuroimaging approach. *Nature Neuroscience*, 19, 1175.
- Gleitman, H., Fridlund, J., and Reisberg, D. (2010). *Psychology* (8th edn.). New York: Norton.
- Glimcher, P. W. (2003). *Decisions, Uncertainty, and the Brain: The Science of Neuroeconomics*. Cambridge MA: MIT Press.
- Glimcher, P. W. (2011) Understanding dopamine and reinforcement learning: The dopamine reward prediction error hypothesis. *Proceedings of the National Academy of Sciences USA*, 108(Suppl 3), 15647–54.
- Glimcher, P. W., and Fehr, E. (2014). *Neuroeconomics*. Cambridge MA: Academic Press.
- Glover, S. R., and Dixon, P. (2001). Dynamic illusion effects in a reaching task: Evidence for separate visual representations in the planning and control of reaching. *Journal of Experimental Psychology: Human Perception and Performance*, 27, 560–72.
- Goense, J., Whittingstall, K., and Logothetis, N. K. (2012). Neural and BOLD responses across the brain. *WIREs Cognitive Science*, 3, 75–86.
- Goldman, A. (2006). *Simulating Minds*. New York: Oxford University Press.
- Goodale, M. A., and Milner, A. D. (2013). *Sight Unseen* (2nd edn.). New York: Oxford University Press.
- Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. Cambridge MA: MIT Press.
- Gopnik, A., and Meltzoff, A. N. (1997). *Words, Thoughts, and Theories*. Cambridge, MA: MIT Press.
- Gordon, R. (1986). Folk psychology as simulation. *Mind and Language*, 1, 158–71.
- Gorman, R. P., and Sejnowski, T. J. (1988). Analysis of hidden units in a layered network trained to identify sonar targets. *Neural Networks*, 1, 75–89.
- Grainger, J., and Jacobs, A. M. (1998). *Localist Connectionist Approaches to Human Cognition*. Mahwah, NJ: Lawrence Erlbaum.

- Graule, M. A., Chirarattananon, P., Fuller, S. B., Jafferis, N. T., Ma, K. Y., Spenko, M., ... Wood, R. J. (2016). Perching and takeoff of a robotic insect on overhangs using switchable electrostatic adhesion. *Science*, 352(6288), 978–82.
- Greenwald, A. G., Draine, S. C., and Abrams, R. L. (1996). Three cognitive markers of unconscious semantic activation. *Science*, 273, 1699–702.
- Griffiths, T. L., Kemp, C., and Tenenbaum, J. B. (2008). Bayesian models of cognition. In R. Son (ed.), *Cambridge Handbook of Computational Cognitive Modeling*. Cambridge: Cambridge University Press.
- Griggs, R. A., and Cox, J. R. (1982). The elusive thematic materials effect in the Wason selection task. *British Journal of Psychology*, 73, 407–20.
- Hacking, I. (2001). *An Introduction to Probability and Inductive Logic*. Cambridge: Cambridge University Press.
- Hadley, R. F. (2000). Cognition and the computational power of connectionist networks. *Connection Science*, 12, 95–110.
- Hampson, R. E., Song, D., Robinson, B. S., Fetterhoff, D., Dakos, A. S., Roeder, B. M., ... Deadwyler, S. A. (2018). Developing a hippocampal neural prosthetic to facilitate human memory encoding and recall. *Journal of Neural Engineering*, 15(3), 036014.
- Hardy-Vallée, B. (2007). Decision-making: A neuroeconomic perspective. *Philosophy Compass*, 2(6), 939–53.
- Harnad, S. (1990). The symbol-grounding problem. *Physica D*, 42, 335–46.
- Haugeland, J. (1985). *Artificial Intelligence: The Very Idea*. Cambridge, MA: MIT Press.
- Haugeland, J. (1997). *Mind Design II: Philosophy, Psychology, Artificial Intelligence*. Cambridge, MA: MIT Press.
- Heal, J. (1986). Replication and functionalism. In J. Butterfield (ed.), *Language, Mind and Logic*. Cambridge: Cambridge University Press.
- Hebb, D. O. (1949). *The Organization of Behavior: A Neuropsychological Theory*. New York: Wiley.
- Heeger, D. J., and Ress, D. (2002). What does fMRI tell us about neuronal activity? *Nature Reviews Neuroscience*, 3, 142–51.
- Heil, J. (2004). *Philosophy of Mind: A Guide and Anthology*. New York: Oxford University Press.
- Henson, R. (2006). Forward inference using functional neuroimaging: Dissociations versus associations. *Trends in Cognitive Sciences*, 10, 64–9.
- Herngenhahn, B. R., and Henley, T. (2013). *An Introduction to the History of Psychology*. Boston: Cengage Learning.
- Hespos, S. J., and van Marle, K. (2012). Physics for infants: Characterizing the origins of knowledge about objects, substances, and number. *WIREs Cognitive Science*, 3, 19–27.
- Heyes, C. (2014). False belief in infancy: A fresh look. *Developmental Science*, 17(5), 647–59.
- Hinton, G. E., McClelland, J. L., and Rumelhart, D. E. (1986). Distributed representations. In D. E. Rumelhart, J. L. McClelland, and the PDP Research Group (eds.), *Parallel Distributed Processing: Explorations in the Microstructures of Cognition: Vol. 1. Foundations*. Cambridge, MA: MIT Press.
- Hirose, S. (1993). *Biologically Inspired Robots: Snake-Like Locomotors and Manipulators*. Oxford: Oxford University Press.
- Hirschfeld, L. A., and Gelman, S. F. (eds.) (1994). *Mapping the Mind: Domain Specificity in Cognition and Culture*. Cambridge: Cambridge University Press.
- Hodges, A. (2014). *Alan Turning: The Enigma*. Princeton, NJ: Princeton University Press.
- Hohwy, J. (2009). The neural correlates of consciousness: New experimental approaches needed? *Consciousness and Cognition*, 18, 428–38.
- Hohwy, J. (2013). *The Predictive Mind*. Oxford: Oxford University Press.

- Hohwy, J., Roepstorff, A., and Friston, K. (2008). Predictive coding explains binocular rivalry: An epistemological review. *Cognition*, 108, 687–701.
- Hopfinger, J. B., Luck, S. J., and Hillyard, S. A. (2004). Selective attention: Electrophysiological and neuromagnetic studies. In M. Gazzaniga (ed.), *The Cognitive Neurosciences* (3rd edn.). Cambridge, MA: MIT Press.
- Horwich, P. (1982). *Probability and Evidence*. Cambridge: Cambridge University Press.
- Houghton, G. (2005). *Connectionist Models in Cognitive Psychology*. Oxford: Oxford University Press.
- Huettel, S. A. (2012). Event-related fMRI in cognition. *Neuroimage*, 62(2), 1152–6.
- Humphreys, G. W., Duncan, J., and Treisman, A. (eds.) (1999). *Attention, Space, and Action: Studies in Cognitive Neuroscience*. Oxford: Oxford University Press.
- Husain, M., and Nachev, P. (2007). Space and the parietal cortex. *Trends in Cognitive Sciences*, 11, 30–6.
- Hutchins, E. (1995). *Cognition in the Wild*. Cambridge, MA: MIT Press.
- Iacoboni, M., and Dapretto, M. (2006). The mirror neuron system and the consequences of its dysfunction. *Nature Reviews Neuroscience*, 7, 942–51.
- Ijspeert, A. J. (2014). Biorobotics: Using robots to emulate and investigate agile locomotion. *Science*, 346(6206), 196–203.
- Isac, D., and Reiss, C. (2013). *I-Language: An Introduction to Linguistics as Cognitive Science* (2nd edn.). Oxford: Oxford University Press.
- Jackson, F. (1982). Epiphenomenal qualia. *Philosophical Quarterly*, 32, 127–36.
- Jackson, F. (1986). What Mary didn't know. *Journal of Philosophy*, 83, 291–5.
- Jackson, F. (2003). Mind and illusion. In A. O'Hear (ed.), *Minds and Persons: Royal Institute of Philosophy Supplement*. Cambridge: Cambridge University Press.
- Jackson, P. (1998). *Introduction to Expert Systems*. Harlow, UK: Addison-Wesley.
- Jacob, P., and Jeannerod, M. (2003). *Ways of Seeing: The Scope and Limits of Visual Cognition*. New York: Oxford University Press.
- Jacobs, R. A., and Kruschke, J. K. (2011). Bayesian learning theory applied to human cognition. *Wiley Interdisciplinary Reviews: Cognitive Science*, 2(1), 8–21.
- Jansen, P. A., and Watter, S. (2012). Strong systematicity through sensorimotor conceptual grounding: An unsupervised, developmental approach to connectionist sentence processing. *Connection Science*, 24(1), 25–55.
- Jeffrey, R. (1983). *The Logic of Decision*. Chicago: University of Chicago Press.
- Jirsa, V. K., and McIntosh, A. R. (eds.) (2007). *The Handbook of Brain Connectivity*. Berlin: Springer.
- Johnson, K. (2004). Gold's theorem and cognitive science. *Philosophy of Science*, 71, 571–92.
- Johnson-Laird, P. N. (1988). *Computer and the Mind: An Introduction to Cognitive Science*. Cambridge, MA: Harvard University Press.
- Jones, J., and Roth, D. (2003). *Robot Programming: A Practical Guide to Behavior-Based Robotics*. New York: McGraw-Hill.
- Kalat, J. W. (2010). *Introduction to Psychology* (9th edn.). Belmont, CA: Wadsworth Thomson Learning.
- Kandel, E. R., Schwarz, J. H., and Jessell, T. M. (2012). *Principles of Neural Science* (5th edn.). New York: McGraw-Hill Medical.
- Kanwisher, N. (2000). Domain specificity in face perception. *Nature Neuroscience*, 3, 759–63.
- Kanwisher, N., McDermott, J., and Chun, M. (1997). The fusiform face area: A module in human extrastriate cortex specialized for the perception of faces. *Journal of Neuroscience*, 17, 4302–11.
- Kaplan, M. (1996). *Decision Theory as Philosophy*. Cambridge: Cambridge University Press.
- Kelly, W. M., Macrae, C. N., Wyland, C. L., Caglar, S., Inati, S., and Heatherton, T. F. (2002). Finding the self? An event-related fMRI study. *Journal of Cognitive Neuroscience*, 14, 785–94.

- Kersten, D., Mamassian, P., and Yuille, A. (2004). Object perception as Bayesian inference. *Annual Review of Psychology*, 55, 271–304.
- Kilner, J. M., and Lemon, R. N. (2013). What we know currently about mirror neurons. *Current Biology*, 23(23), R1057–62.
- Kiran, S., and Lebel, K. (2007). Crosslinguistic semantic and translation priming in normal bilingual individuals and bilingual aphasia. *Clinical Linguistics and Phonetics*, 4, 277–303.
- Knill, D. C., and Whitman, R. (2008). *Perception as Bayesian Inference*. Cambridge: Cambridge University Press.
- Koch, C. (2004). *The Quest for Consciousness: A Neurobiological Approach*. Englewood, CO: Roberts.
- Koch, C., and Tsuchiya, N. (2007). Attention and consciousness: Two distinct brain processes. *Trends in Cognitive Sciences*, 11, 229–35.
- Koch, C., Massimini, M., Boly, M., and Tononi, G. (2016). Neural correlates of consciousness: Progress and problems. *Nature Reviews Neuroscience*, 17, 307.
- Koh, J.-S., Yang, E., Jung, G.-P., Jung, S.-P., Son, J. H., Lee, S.-I., ... Cho, K.-J. (2015). Jumping on water: Surface tension-dominated jumping of water striders and robotic insects. *Science*, 349(6247), 517–21.
- Kording, K. P., and Wolpert, D. M. (2004). Bayesian integration in sensorimotor learning. *Nature*, 427, 244.
- Kording, K. P., and Wolpert, D. M. (2006). Bayesian decision theory in sensorimotor control. *Trends in Cognitive Sciences*, 10(7), 319–26.
- Kosslyn, S. M. (1973). Scanning visual images: Some structural implications. *Perception and Psychophysics*, 14, 341–70.
- Kosslyn, S. M., Thompson, W. L., and Ganis, G. (2006). *The Case for Mental Imagery*. Oxford: Oxford University Press.
- Kotz, S. A. (2001). Neurolinguistic evidence for bilingual language representation: A comparison of reaction times and event-related brain potentials. *Bilingualism: Language and Cognition*, 4, 143–54.
- Kouider, S., and Dehaene, S. (2007). Levels of processing during non-conscious perception: A critical review of visual masking. *Philosophical Transactions of the Royal Society of London B*, 362(1481), 857–75.
- Kouider, S., de Gardelle, V., Sackur, J., and Dupoux, E. (2010). How rich is consciousness? The partial awareness hypothesis. *Trends in Cognitive Sciences*, 14, 301–7.
- Kouider, S., Dehaene, S., Jobert, A., and Le Bihan, D. (2007). Cerebral bases of subliminal and supraliminal priming during reading. *Cerebral Cortex*, 17, 2019–29.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). *Image classification with deep convolutional neural networks*. Paper presented at Advances in Neural Information Processing Systems (NIPS 2012).
- Laird, J. E. (2012). *The Soar Cognitive Architecture*. Cambridge, MA: MIT Press.
- Lake, B., Salakhutdinov, R., and Tenenbaum, J. B. (2016). Human level concept learning through probabilistic program induction. *Science*, 350, 1332–8.
- Lamme, V. A. F. (2003). Why visual attention and awareness are different. *Trends in Cognitive Sciences*, 7, 12–18.
- Lamme, V. A. F. (2006). Towards a true neural stance on consciousness. *Trends in Cognitive Sciences*, 10, 494–501.
- Lashley, K. S. (1951). The problem of serial order in behavior. In A. L. Jeffress (ed.), *Cerebral Mechanisms in Behavior: The Hixon Symposium*. New York: Wiley.
- Laureys, S. (2005). The neural correlate of (un)awareness: Lessons from the vegetative state. *Trends in Cognitive Sciences*, 9, 556–9.

- Lavie, N. (2005). Distracted and confused? Selective attention under load. *Trends in Cognitive Sciences*, 9, 75–82.
- Lazarou, I., Nikolopoulos, S., Petrantonakis, P. C., Kompatsiaris, I., and Tsolaki, M. (2018). EEG-based brain-computer interfaces for communication and rehabilitation of people with motor impairment: A novel approach of the 21st century. *Frontiers in Human Neuroscience*, 12(14).
- Lebiere, C. (2003). ACT. In L. Nadel (ed.), *Encyclopedia of Cognitive Science*. New York: Nature.
- LeCun, Y., Bengio, Y., and Hinton, G. (2015). Deep learning. *Nature*, 521, 436.
- Leslie, A. M. (1987). Pretense and representation: The origins of “theory of mind.” *Psychological Review*, 94, 412–26.
- Leslie, A. M., and Polizzi, P. (1998). Inhibitory processing in the false belief task: Two conjectures. *Developmental Science*, 1, 247–53.
- Leslie, A. M., Friedman, O., and German, T. P. (2004). Core mechanisms in “theory of mind.” *Trends in Cognitive Sciences*, 8, 529–33.
- Leslie, A. M., German, T. P., and Polizzi, P. (2005). Belief-desire reasoning as a process of selection. *Cognitive Psychology*, 50, 45–85.
- Levine, J. (1983). Materialism and qualia: The explanatory gap. *Pacific Philosophical Quarterly*, 64, 354–61.
- Libby, T., Moore, T. Y., Chang-Siu, E., Li, D., Cohen, D. J., Jusufi, A., and Full, R. J. (2012). Tail-assisted pitch control in lizards, robots and dinosaurs. *Nature*, 481(7380), 181–4.
- Lidz, J., Waxman, S., and Freedman, J. (2003). What infants know about syntax but couldn’t have learned: Experimental evidence for syntactic structure at 18 months. *Cognition*, 89(3), B65–73.
- Logothetis, N. K. (2001). The underpinnings of the BOLD functional magnetic resonance imaging signal. *Journal of Neuroscience*, 23, 3963–71.
- Logothetis, N. K. (2002). The neural basis of the blood-oxygen-level-dependent functional magnetic resonance imaging signal. *Philosophical Transactions of the Royal Society of London B*, 357(1424), 1003–37.
- Logothetis, N. K. (2008). What we can do and what we cannot do with fMRI. *Nature*, 453, 869–78.
- Logothetis, N. K., Pauls, J., Augath, M., Trinath, T., and Oeltermann, A. (2001). Neurophysiological investigation of the fMRI signal. *Nature*, 412, 150–7.
- Lovett, M. C., and Anderson, J. R. (2005). Thinking as a production system. In K. J. Holyoak and R. G. Morrison (eds.), *The Cambridge Handbook of Thinking and Reasoning*. Cambridge: Cambridge University Press.
- Low, J., and Perner, J. (2012). Implicit and explicit theory of mind: State of the art. *British Journal of Developmental Psychology*, 30, 1–30.
- Luck, S. J. (2005). *An Introduction to the Event-Related Potential Technique*. Cambridge, MA: MIT Press.
- Luck, S. J., and Ford, M. A. (1998). On the role of selective attention in visual perception. *Proceedings of the National Academy of Sciences USA*, 95, 825–30.
- Luck, S. J., and Kappenman, E. S. (2011). *The Oxford Handbook of Event-Related Potential Components*. Oxford: Oxford University Press.
- Ludlow, P., Nagasawa, Y., and Stoljar, D. (eds.) (2004). *There’s Something about Mary*. Cambridge, MA: MIT Press.
- Luo, Y., and Baillargeon, R. (2010). Toward a mentalistic account of early psychological reasoning. *Current Directions in Psychological Science*, 19, 301–7.
- Luria, A. R. (1970). The functional organization of the brain. *Scientific American*, 222, 66–72.
- Macdonald, C., and Macdonald, G. (1995). *Connectionism*. Oxford: Blackwell.
- Machery, E. (2008). Massive modularity and the flexibility of human cognition. *Mind and Language*, 23, 263–72.

- Machery, E. (2012). Dissociations in neuropsychology and cognitive neuroscience. *Philosophy of Science*, 79, 490–518.
- Mack, A., and Rock, I. (1998). *Inattentional Blindness*. Cambridge, MA: MIT Press.
- Manassi, M., Sayim, B., and Herzog, M. H. (2013). When crowding of crowding leads to uncrowding. *Journal of Vision*, 13(13), 10. doi:10.1167/13.13.10.
- Marcus, G. (2003). *The Algebraic Mind: Integrating Connectionism and Cognitive Science*. Cambridge, MA: MIT Press.
- Marcus, G., Ullman, M., Pinker, S., Hollander, M., Rosen, T. J., and Xu, F. (1992). *Overregularization in Language Acquisition*. Chicago: University of Chicago Press.
- Mareschal, D., and Johnson, S. P. (2002). Learning to perceive object unity: A connectionist account. *Developmental Science*, 5, 151–85.
- Mareschal, D., Plunkett, K., and Harris, P. (1995). Developing object permanence: A connectionist model. In J. D. Moore and J. F. Lehman (eds.), *Proceedings of the Seventeenth Annual Conference of the Cognitive Science Society*. Mahwah, NJ: Lawrence Erlbaum.
- Margolis, E., Samuels, R., and Stich, S. (eds.) (2012). *The Oxford Handbook of Philosophy of Cognitive Science*. Oxford: Oxford University Press.
- Marr, D. (1982). *Vision: A Computational Investigation into the Human Representation and Processing of Visual Information*. San Francisco: W. H. Freeman.
- Marr, D. (2010). *Vision: A Computational Investigation into the Human Representation and Processing of Visual Information*. 1982; repr., Cambridge, MA: MIT Press.
- Marr, D., and Hildreth, E. (1980). Theory of edge detection. *Proceedings of the Royal Society of London*, 204, 187–217.
- Marshall, J. C., and Halligan, P. W. (1988). Blindsight and insight in visuospatial neglect. *Nature*, 366, 766–7.
- Martens, S., and Wyble, B. (2010). The attentional blink: Past, present, and future of a blind spot in perceptual awareness. *Neuroscience and Biobehavioral Reviews*, 34, 947–57.
- Matarić, M. (1997). Behavior-based control: Examples from navigation, learning, and group behavior. *Journal of Experimental and Theoretical Artificial Intelligence*, 9, 323–36.
- Matarić, M. (1998). Behavior-based robotics as a tool for synthesis of artificial behavior and analysis of natural behavior. *Trends in Cognitive Sciences*, 2, 82–7.
- Matarić, M. (2007). *The Robotics Primer*. Cambridge, MA: MIT Press.
- Mathersul, D., McDonald, S., and Rushby, J. A. (2013). Understanding advanced theory of mind and empathy in high-functioning adults with autism spectrum disorder. *Journal of Clinical and Experimental Neuropsychology*, 35(6), 655–68.
- McClelland, J. L., and Jenkins, E. (1991). Nature, nurture, and connectionism: Implications for connectionist models of development. In K. van Lehn (ed.), *Architectures for Intelligence: The 22nd (1988) Carnegie Symposium on Cognition*. Hillsdale, NJ: Lawrence Erlbaum.
- McClelland, J. L., and Patterson, K. (2002). Rules or connections in past-tense inflections: What does the evidence rule out? *Trends in Cognitive Sciences*, 6, 465–72.
- McClelland, J. L., Rumelhart, D. E., and the PDP Research Group (1986). *Parallel Distributed Processing: Explorations in the Microstructures of Cognition. Vol. 2: Psychological and Biological Models*. Cambridge, MA: MIT Press.
- McClelland, J. L., Botvinick, M. M., Noelle, D. C., et al. (2010). Letting structure emerge: Connectionist and dynamical systems approaches to cognition. *Trends in Cognitive Sciences*, 14, 348–56.
- McCulloch, W. S., and Pitts, W. H. (1943). A logical calculus of the ideas immanent in nervous activity. *Bulletin of Mathematical Biophysics*, 5, 115–33.
- McDermott, J. H. (2009). The cocktail party problem. *Current Biology*, 19, R1024–7.

- McLeod, P., Plunkett, K., and Rolls, E. T. (1998). *Introduction to the Connectionist Modelling of Cognitive Processes*. Oxford: Oxford University Press.
- McSweeney, F. K., and Murphy, E. S. (2014). *The Wiley-Blackwell Handbook of Classical and Operant Conditioning*. Chichester, UK: Wiley-Blackwell.
- Medsker, L. R., and Schulte, T. W. (2003). Expert systems. In L. Nadel (ed.), *Encyclopedia of Cognitive Science* (vol. 2). New York: Nature.
- Melcher, D., and Colby, C. L. (2008). Trans-saccadic perception. *Trends in Cognitive Sciences*, 12, 466–73.
- Merikle, P. M., Joordens, S., and Stoltz, J. (1995). Measuring the relative magnitude of unconscious influences. *Consciousness and Cognition*, 4, 422–39.
- Metzinger, T. (ed.) (2000). *Neural Correlates of Consciousness: Empirical and Conceptual Issues*. Cambridge, MA: MIT Press.
- Michalski, R. S., and Chilausky, R. L. (1980). Learning by being told and learning from examples: An experimental comparison of the two methods for knowledge acquisition in the context of developing an expert system for soybean disease diagnosis. *International Journal of Policy Analysis and Information Systems*, 4, 125–61.
- Miller, G. A. (1956). The magical number seven, plus or minus two: Some limits on our capacity for processing information. *Psychological Review*, 63, 81–97.
- Miller, G. A. (2003). The cognitive revolution: A historical perspective. *Trends in Cognitive Sciences*, 7, 141–4.
- Millican, P., and Clark, A. (1996). *The Legacy of Alan Turing: Machines and Thought*. London: Clarendon Press.
- Milner, A. D. (2012). Is visual processing in the dorsal stream accessible to consciousness? *Proceedings of the Royal Society B*, 279, 2289–98.
- Milner, A. D., and Goodale, M. A. (1998). *The Visual Brain in Action* (Precis). *Psyche*, 4.
- Milner, A. D., and Goodale, M. A. (2006). *The Visual Brain in Action* (2nd edn.). Oxford: Oxford University Press.
- Milner, A. D., and Goodale, M. A. (2008). Two visual systems reviewed. *Neuropsychologia*, 46, 774–85.
- Milner, B. (1966). Amnesia following operation on the temporal lobes. In C. W. M. Whitty and O. L. Zangwill (eds.), *Amnesia*. London: Butterworth.
- Minati, L., Varotto, G., D'Incerti, L., Panzica, F., and Chan, D. (2013). From brain topography to brain topology: Relevance of graph theory to functional neuroscience. *Neuroreport*, 24(10), 536–43.
- Minsky, M., and Papert, S. (1969). *Perceptrons*. Cambridge, MA: MIT Press.
- Mishkin, M. L., Ungerleider, G., and Macko, K. A. (1983). Object vision and spatial vision: Two cortical pathways. *Trends in Neurosciences*, 6, 414–17. Reprinted in W. Bechtel, P. Mandik, J. Mundale, and R. Stufflebeam (eds.) (2001), *Philosophy and the Neurosciences: A Reader*. Oxford: Blackwell.
- Mitchell, J. P., Banaji, M. R., and Macrae, C. N. (2005). The link between social cognition and self-referential thought in the medial prefrontal cortex. *Journal of Cognitive Neuroscience*, 17, 1306–15.
- Mitchell, T. M. (1997). *Machine Learning*. Boston: McGraw-Hill.
- Mohan, A., Roberto, A. J., Mohan, A., Lorenzo, A., Jones, K., Carney, M. J., ... Lapidus, K. A. (2016). The significance of the default mode network (DMN) in neurological and neuropsychiatric disorders: A review. *Yale Journal of Biology and Medicine*, 89(1), 49–57.
- Molenberghs, P., Sale, M. V., and Mattingley, J. B. (2012). Is there a critical lesion site for unilateral spatial neglect? A meta-analysis using activation likelihood estimation. *Frontiers in Human Neuroscience*, 6, 1–10.
- Montague, P., Dayan, P., and Sejnowski, T. (1996). A framework for mesencephalic dopamine systems based on predictive Hebbian learning. *Journal of Neuroscience*, 16(5), 1936–47.

- Morse, S. J., and Roskies, A. L. (2013). *A Primer on Criminal Law and Neuroscience*. Oxford: Oxford University Press.
- Mountcastle, V. B., Lynch, J. C., Georgopoulos, A., Sakata, H., and Acuna, C. (1975). Posterior parietal association cortex of the monkey: Command functions for operations within extrapersonal space. *Journal of Neurophysiology*, 38(4), 871–908.
- Mukamel, R., Gelbard, H., Arieli, A., Hasson, U., Fried, I., and Malach, R. (2005). Coupling between neuronal firing, field potentials, and fMRI in human auditory cortex. *Science*, 309, 951–4.
- Munakata, Y. (2001). Graded representations in behavioral dissociations. *Trends in Cognitive Sciences*, 5, 309–15.
- Munakata, Y., and McClelland, J. L. (2003). Connectionist models of development. *Developmental Science*, 6, 413–29.
- Munakata, Y., McClelland, J. L., Johnson, M. H., and Siegler, R. S. (1997). Rethinking infant knowledge: Toward an adaptive process account of successes and failures in object permanence tasks. *Psychological Review*, 104, 686–713.
- Nadel, L. (ed.) (2005). *Encyclopedia of Cognitive Science*. Chichester, UK: Wiley.
- Needham, A., and Libertus, K. (2011). Embodiment in early development. *WIREs Cognitive Science*, 2, 117–23.
- Newmeyer, F. J. (1986). *Linguistic Theory in America*. London: Academic Press.
- Nichols, S., Stich, S., Leslie, A., and Klein, D. (1996). Varieties of off-line simulation. In P. Carruthers and P. K. Smith (eds.), *Theories of Theory of Mind*. Cambridge: Cambridge University Press.
- Nilsson, N. J. (1984). Shakey the robot. SRI International, Technical Note 323.
- Norman, D. A., and Shallice, T. (1980). Attention to action: Willed and automatic control of behaviour. Reprinted in M. Gazzaniga (ed.), *Cognitive Neuroscience: A Reader*. Oxford: Blackwell (2000).
- Oakes, L. M. (2010). Using habituation of looking time to assess mental processes in infancy. *Journal of Cognition and Development*, 11, 255–68.
- Oaksford, M., and Chater, N. (1994). A rational analysis of the selection task as optimal data selection. *Psychological Review*, 101, 608–31.
- Oberauer, K. (2006). Reasoning with conditionals: A test of formal models of four theories. *Cognitive Psychology*, 53, 238–83.
- O'Grady, W., Archibald, J., Aronoff, M., and Rees-Miller, J. (2010). *Contemporary Linguistics: An Introduction* (6th edn.). Boston: Bedford/St. Martin's.
- Onishi, K. H., and Baillargeon, R. (2005). Do 15-month-old infants understand false beliefs? *Science*, 308, 255–8.
- Orban, G. A., Van Essen, D., and Vanduffel, W. (2004). Comparative mapping of higher visual areas in monkeys and humans. *Trends in Cognitive Sciences*, 8, 315–24.
- O'Reilly, R. C., and Munakata, Y. (2000). *Computational Explorations in Computational Neuroscience: Understanding the Mind by Simulating the Brain*. Cambridge, MA: MIT Press.
- Owen, A. M., Coleman, M. R., Boly, M., Davis, M. H., Laureys, S., and Pickard, J. D. (2006). Detecting awareness in the vegetative state. *Science*, 313, 1402.
- Page, M. (2000). Connectionist modeling in psychology: A localist manifesto. *Behavioral and Brain Sciences*, 23, 443–67.
- Passingham, R. (2009). How good is the macaque monkey model of the human brain? *Current Opinion in Neurobiology*, 19, 6–11.
- Pearl, L., and Goldwater, S. (2016). Statistical learning, inductive bias, and Bayesian inference in language acquisition. In J. Lidz, W. Snyder, & C. Pater (eds.), *The Oxford Handbook of Developmental Linguistics*. Oxford: Oxford University Press.
- Pelucchi, B., Hay, J. F., and Saffran, J. R. (2009). Statistical learning in a natural language by 8-month-old infants. *Child Development*, 80(3), 674–85.

- Perner, J. (1993). *Understanding the Representational Mind*. Cambridge, MA: MIT Press.
- Perner, J., and Leekam, S. (2008). The curious incident of the photo that was accused of being false: Issues of domain specificity in development, autism, and brain imaging. *Quarterly Journal of Experimental Psychology*, 61, 76–89.
- Perner, J., and Roessler, J. (2012). From infants' to children's appreciation of belief. *Trends in Cognitive Sciences*, 16, 519–25.
- Perner, J., and Ruffman, T. (2005). Infants' insight into the mind: How deep? *Science*, 308(5719), 214–16.
- Perner, J., Leekam, S. R., and Wimmer, H. (1987). Three-year-olds' difficulty with false belief: The case for a conceptual deficit. *British Journal of Developmental Psychology*, 5(2), 125–37.
- Perone, S., and Simmering, V. R. (2017). Applications of dynamic systems theory to cognition and development: New frontiers. *Advances in Child Development and Behavior*, 52, 43–80.
- Peru, A., Moro, V., Avesani, R., and Aglioti, S. (1996). Overt and covert processing of left-side information in unilateral neglect investigated with chimeric drawings. *Journal of Clinical and Experimental Neuropsychology*, 18, 621–30.
- Petersen, S. E., and Fiez, J. A. (2001). The processing of single words studied with positron emission tomography. In W. Bechtel, P. Mandik, J. Mundale, and R. S. Stufflebeam (eds.), *Philosophy and the Neurosciences: A Reader*. Malden, MA: Blackwell.
- Petersen, S. E., Fox, P. T., Posner, M. I., and Mintun, M. (1988). Positron emission tomographic studies of the cortical anatomy of single-word processing. *Nature*, 331, 585–9.
- Petzold, C. (2008). *The Annotated Turing: A Guided Tour through Alan Turing's Historic Paper on Computability and the Turing Machine*. Indianapolis, IN: Wiley.
- Pfeifer, R., Iida, F., and Gómez, G. (2006). Morphological computation for adaptive behavior and cognition. *International Congress Series*, 1291, 22–9.
- Phillips, M. L., Young, A. W., Senior, C., et al. (1997). A specific neural substrate for perceiving facial expressions of disgust. *Nature*, 389, 495–8.
- Piaget, J. (1954). *The Construction of Reality in the Child*. New York: Basic Books.
- Piccinini, G. (2004). The first computational theory of mind and brain: A close look at McCulloch and Pitts' "Logical calculus of the ideas immanent in nervous activity." *Synthese*, 141, 175–215.
- Piccinini, G., and Craver, C. (2011). Integrating psychology and neuroscience: Functional analyses as mechanism sketches. *Synthese*, 183, 283–311.
- Pinker, S. (1994). *The Language Instinct: How the Mind Creates Language*. New York: Harper.
- Pinker, S. (1997). *How the Mind Works*. New York: Norton.
- Pinker, S. (2005). So how does the mind work? *Mind and Language*, 20, 1–24.
- Pinker, S., and Prince, A. (1988a). On language and connectionism: Analysis of a parallel distributed processing model of language acquisition. *Cognition*, 28, 73–193.
- Pinker, S., and Prince, A. (1988b). Rules and connections in human language. In R. Morris (ed.), *Parallel Distributed Processing*. Oxford: Oxford University Press.
- Pinker, S., and Ullman, M. T. (2002). The past and future of the past tense. *Trends in Cognitive Sciences*, 6, 456–63.
- Platt, M. L., and Glimcher, P. W. (1999). Neural correlates of decision variables in parietal cortex. *Nature*, 400(6741), 233–8.
- Plaut, D. C., and McClelland, J. L. (2010). Locating object knowledge in the brain: Comment on Bowers's (2009) attempt to revive the grandmother cell hypothesis. *Psychological Review*, 117, 284–90.
- Plaut, D. C., Banich, M. T., and Mack, M. (2003). Connectionist modeling of language: Examples and implications. In M. T. Banich and M. Mack (eds.), *Mind, Brain, and Language: Multidisciplinary Perspectives*. Mahwah, NJ: Lawrence Erlbaum.

- Plotnik, R., and Kouyoumdjian, H. (2010). *Introduction to Psychology* (9th edn.). Belmont, CA: Wadsworth Thomson Learning.
- Plunkett, K., and Elman, J. L. (1997). *Exercises in Rethinking Innateness: A Handbook for Connectionist Simulations*. Cambridge, MA: MIT Press.
- Plunkett, K., and Marchman, V. (1993). From rote learning to system building: Acquiring verb morphology in children and connectionist nets. *Cognition*, 48, 21–69.
- Poldrack, R. A. (2006). Can cognitive processes be inferred from neuroimaging data? *Trends in Cognitive Sciences*, 10, 59–63.
- Poldrack, R. A. (2018). *The New Mind Readers: What Neuroimaging Can and Cannot Reveal about Our Thoughts*. Princeton, NJ: Princeton University Press.
- Poldrack, R. A., Mumford, J. A., and Nichols, T. E. (2011). *Handbook of Functional MRI Data Analysis*. Cambridge: Cambridge University Press.
- Pollard, P., and Evans, J. St. B. T. (1987). Content and context effects in reasoning. *American Journal of Psychology*, 100, 41–60.
- Poole, D. L., and Mackworth, A. K. (2010). *Artificial Intelligence: Foundations of Computational Agents*. Cambridge: Cambridge University Press.
- Pöppel, E., Frost, D., and Held, R. (1973). Residual visual function after brain wounds involving the central visual pathways in man. *Nature*, 243, 295–6.
- Port, R. F., and Van Gelder, T. (1995). *Mind as Motion: Explorations in the Dynamics of Cognition*. Cambridge, MA: MIT Press.
- Posner, M. I. (1980). Orienting of attention. *Quarterly Journal of Experimental Psychology*, 32, 3–25.
- Posner, M. I. (1989). *Foundations of Cognitive Science*. Cambridge, MA: MIT Press.
- Posner M. I. (ed.) (2004). *The Cognitive Neuroscience of Attention*. New York: Guilford.
- Posner, M. I. (2017). Attentional mechanisms. Reference Module in Neuroscience and Biobehavioral Psychology: Elsevier. [online resource]
- Posner, M. I., and Raichle, M. E. (1994). *Images of Mind*. New York: Scientific American Library.
- Poulin-Dubois, D., Brooker, I., and Chow, V. (2009). The developmental origins of naive psychology in infancy. *Advances in Child Development and Behavior*, 37, 55–104.
- Preston, J. M., and Bishop, J. M. (2002). *Views into the Chinese Room: New Essays on Searle and Artificial Intelligence*. Oxford: Oxford University Press.
- Prinz, J. (2012). *The Conscious Brain*. New York: Oxford University Press.
- Pullum Geoffrey, K., and Scholz Barbara, C. (2002). Empirical assessment of stimulus poverty arguments. *The Linguistic Review*, 18, 9–50.
- Purves, D., Augustine, G. J., Fitzpatrick, D., Hall, W. C., Anthony-Samuel, L., and White, L. E. (2011) *Neuroscience* (5th edn.). Sunderland, MA: Sinauer Associates.
- Pylyshyn, Z. (1980). Computation and cognition: Issues in the foundations of cognitive science. *Behavioral and Brain Sciences*, 3, 111–69.
- Pylyshyn, Z. (1984). *Computation and Cognition: Toward a Foundation for Cognitive Science*. Cambridge, MA: MIT Press.
- Pylyshyn, Z. (ed.) (1987). *The Robot's Dilemma: The Frame Problem in Artificial Intelligence*. Norwood, NJ: Ablex.
- Quinlan, P. T., van der Maas, H. L. J., Jansen, B. R. J., Booij, O., and Rendell, M. (2007). Re-thinking stages of cognitive development: An appraisal of connectionist models of the balance scale task. *Cognition*, 103, 413–59.
- Raffone, A., Srinivasan, N., and van Leeuwen, C. (2014). The interplay of attention and consciousness in visual search, attentional blink and working memory consolidation. *Philosophical Transactions of the Royal Society B*, 369(1641), 20130215.

- Raichle, M. E., and Mintun, M. A. (2006). Brain work and brain imaging. *Annual Review of Neuroscience*, 29, 449–76.
- Rakoczy, H. (2012). Do infants have a theory of mind? *British Journal of Development Psychology*, 30(Pt 1), 59–74.
- Ramnani, N., Behrens, T. E. J., Penny, W., and Matthews, P. M. (2004). New approaches for exploring functional and anatomical connectivity in the human brain. *Biological Psychiatry*, 56, 613–19.
- Ramsey, W., Stich, S. P., and Rumelhart, D. E. (1991). *Philosophy and Connectionist Theory*. Hillsdale, NJ: Lawrence Erlbaum.
- Rees, G., Friston, K., and Koch, C. (2000). A direct quantitative relationship between the functional properties of human and macaque V5. *Nature Neuroscience*, 3, 716–23.
- Regier, T., and Gahl, S. (2004). Learning the unlearnable: The role of missing evidence. *Cognition*, 93(2), 147–55.
- Riley, M. A., and Holden, J. G. (2012). Dynamics of cognition. *WIREs Cognitive Science*, 3, 593–606.
- Ritter, F. E. (2003). Soar. In L. Nadel (ed.), *Encyclopedia of Cognitive Science*. New York: Nature.
- Rizzolatti, G., and Fogassi, L. (2014). The mirror mechanism: Recent findings and perspectives. *Philosophical Transactions of the Royal Society of London, Series*, 369(1644), 20130420.
- Rizzolatti, G., and Sinigaglia, C. (2008). *Mirrors in the Brain: How Our Minds Share Actions and Emotions*. Trans. F. Anderson. Oxford: Oxford University Press.
- Rizzolatti, G., and Sinigaglia, C. (2010). The functional role of the parieto-frontal mirror circuit: Interpretations and misinterpretations. *Nature Reviews Neuroscience*, 11, 264–74.
- Rizzolatti, G., Fogassi, L., and Gallese, V. (2001). Neurophysiological mechanisms underlying the understanding and imitation of action. *Nature Reviews Neuroscience*, 2, 661–70.
- Rizzolatti, G., Fogassi, L., and Gallese, V. (2006). Mirrors of the mind. *Scientific American*, 295, 54–61.
- Robbins, R., and Aydede, M. (eds.) (2008). *The Cambridge Handbook of Situated Cognition*. Cambridge: Cambridge University Press.
- Robinson, D. L., Goldberg, M. E., and Stanton, G. B. (1978). Parietal association cortex in the primate: Sensory mechanisms and behavioral modulations. *Journal of Neurophysiology*, 41(4), 910–32.
- Roediger, H. L., Dudai, Y., and Fitzpatrick, S. M. (2007). *Science of Memory: Concepts*. Oxford: Oxford University Press.
- Roeyers, H., and Demurie, E. (2010). How impaired is mind-reading in high-functioning adolescents and adults with autism? *European Journal of Developmental Psychology*, 7(1), 123–34.
- Rogers, R. (1971). *Mathematical Logic and Formalized Theories*. Amsterdam: North-Holland.
- Rogers, T. T., and McClelland, J. L. (2004). *Semantic Cognition: A Parallel Distributed Processing Approach*. Cambridge, MA: MIT Press.
- Rohde, D., and Plaut, D. C. (1999). Language acquisition in the absence of explicit negative evidence: How important is starting small? *Cognition*, 72, 67–109.
- Rollins, M. (1989). *Mental Imagery: The Limits of Cognitive Science*. Cambridge, MA: MIT Press.
- Rolls, E. T., and Milward, T. (2000). A model of invariant object recognition in the visual system: Learning rules, activation functions, lateral inhibition, and information-based performance measures. *Neural Computation*, 12, 2547–72.
- Rosenblatt, F. (1958). The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65, 386–408.
- Rösler, F., Ranganath, C., Röder, B., and Kluwe, R. (2009). *Neuroimaging of Human Memory: Linking Cognitive Processes to Neural Systems*. New York: Oxford University Press.

- Rossetti, Y., Pisella, L., and McIntosh, R. D. (2017). Rise and fall of the two visual systems theory. *Annals of Physical and Rehabilitation Medicine*, 60(3), 130–40.
- Rowe, J. B., and Frackowiak, R. S. J. (2003). Neuroimaging. In L. Nadel (ed.), *Encyclopedia of Cognitive Science*. New York: Nature.
- Rumelhart, D. E. (1989). The architecture of mind: A connectionist approach. In M. I. Posner (ed.), *Foundations of Cognitive Science*. Cambridge, MA: MIT Press. Reprinted in J. Haugeland (ed.) (1997), *Mind Design II: Philosophy, Psychology, Artificial Intelligence*. Cambridge, MA: MIT Press.
- Rumelhart, D. E., and McClelland, J. L. (1986). On learning the past tenses of English verbs. In J. L. McClelland, D. E. Rumelhart, and the PDP Research Group (eds.), *Parallel Distributed Processing: Explorations in the Microstructures of Cognition: Vol. 2. Psychological and Biological Models*. Cambridge, MA: MIT Press.
- Rumelhart, D. E., McClelland, J. L., and the PDP Research Group (eds.) (1986). *Parallel Distributed Processing: Explorations in the Microstructures of Cognition: Vol. 1. Foundations*. Cambridge, MA: MIT Press.
- Russell, S. J., and Norvig, P. (2003). *Artificial Intelligence: A Modern Approach* (2nd edn.). Upper Saddle River, NJ: Prentice Hall.
- Russell, S. J., and Norvig, P. (2009). *Artificial Intelligence: A Modern Approach* (3rd edn.). New Delhi: Prentice Hall.
- Saffran, J. R., Aslin, R. N., and Newport, E. L. (1996). Statistical learning by 8-month-old infants. *Science*, 274(5294), 1926–8.
- Samson, D., Apperly, I. A., Chiavarino, C., and Humphreys, G. W. (2004). Left temporoparietal junction is necessary for representing someone else's belief. *Nature Neuroscience*, 7, 499–500.
- Samson, D., Apperly, I. A., Kathirgamanathan, U., and Humphreys, G. W. (2005). Seeing it my way: A case of a selective deficit in inhibiting self-perspective. *Brain: A Journal of Neurology*, 128, 1102–11.
- Samuelson, L. K., Jenkins, G. W., and Spencer, J. P. (2015). Grounding cognitive-level processes in behavior: The view from dynamic systems theory. *Topics in Cognitive Science*, 7, 191–205.
- Saxe, R. (2009). Theory of mind (neural basis). In W. Banks (ed.), *Encyclopedia of Consciousness*. Cambridge, MA: MIT Press.
- Saxe, R., and Kanwisher, N. (2005). People thinking about thinking people: The role of the temporo-parietal junction in "Theory of Mind." In J. T. Cacioppo and G. G. Berntson (eds.), *Social Neuroscience: Key Readings*. New York: Psychology Press.
- Saxe, R., Carey, S., and Kanwisher, N. (2004). Understanding other minds: Linking developmental psychology and functional neuroimaging. *Annual Review of Psychology*, 55, 87–124.
- Schenk, T., and McIntosh, R. D. (2010). Do we have independent visual streams for perception and action? *Cognitive Neuroscience*, 1, 52–78.
- Schlatter, M., and Aizawa, K. (2008). Walter Pitts and "A logical calculus." *Synthese*, 162, 235–50.
- Schneider, S. (2011). *The Language of Thought: A New Philosophical Direction*. Cambridge, MA: MIT Press.
- Schneider, S., and Katz, M. (2012). Rethinking the language of thought. *WIREs Cognitive Science*, 3, 153–62.
- Schoonbaert, S., Duyck, W., Brysbaert, M., and Hartsuiker, R. J. (2009). Semantic and translation priming from a first language to a second and back: Making sense of the findings. *Memory & Cognition*, 37, 569–86.
- Schurz, M., and Perner, J. (2015). An evaluation of neurocognitive models of theory of mind. *Frontiers in Psychology*, 6, 1610.
- Schyns, P. G., Gosselin, F., and Smith, M. L. (2008). Information processing algorithms in the brain. *Trends in Cognitive Sciences*, 13, 20–6.

- Scott, R. M., and Baillargeon, R. (2017). Early false-belief understanding. *Trends in Cognitive Sciences*, 21(4), 237–49.
- Searle, J. (1980). Minds, brains, and programs. *Behavioral and Brain Sciences*, 3, 417–57.
- Searle, J. (2004). *Mind: A Brief Introduction*. New York: Oxford University Press.
- Setoh, P., Scott, R. M., and Baillargeon, R. (2016). Two-and-a-half-year-olds succeed at a traditional false-belief task with reduced processing demands. *Proceedings of the National Academy of Sciences USA*, 113(47), 13360–5.
- Shadmehr, R., and Krakauer, J. W. (2008). A computational neuroanatomy for motor control. *Experimental Brain Research*, 185, 359–81.
- Shallice, T., and Warrington, E. K. (1970). Independent functioning of memory stores: A neuropsychological study. *Quarterly Journal of Experimental Psychology*, 22, 261–73.
- Shanahan, M. P. (2003). The frame problem. In L. Nadel (ed.), *Encyclopedia of Cognitive Science*. New York: Nature.
- Shannon, C. E. (1948). A mathematical theory of communication. *Bell System Technical Journal*, 27, 379–423, 623–56.
- Shapiro, L. (2007). The embodied cognition research programme. *Philosophy Compass*, 2, 338–46.
- Shapiro, L. (2011). *Embodied Cognition*. New York: Routledge.
- Shepard, R. N., and Metzler, J. (1971). Mental rotation of three-dimensional objects. *Science*, 171, 701–3.
- Shepherd, G. (1994). *Neurobiology* (3rd edn.). New York: Oxford University Press.
- Siegelmann, H., and Sontag, E. (1991). Turing computability with neural nets. *Applied Mathematics Letters*, 4, 77–80.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., van den Driessche, G., ... Hassabis, D. (2016). Mastering the game of Go with deep neural networks and tree search. *Nature*, 529, 484.
- Silver, D., Schrittwieser, J., Simonyan, K., Antonoglou, I., Huang, A., Guez, A., ... Hassabis, D. (2017). Mastering the game of Go without human knowledge. *Nature*, 550, 354.
- Simons, D., and Chabris, C. (1999). Gorillas in our midst: Sustained inattentional blindness for dynamic events. *Perception*, 28, 1059–74.
- Simons, D., and Rensink, R. A. (2005). Change blindness: Past, present, and future. *Trends in Cognitive Sciences*, 9, 16–20.
- Singer, W. (1999). Neuronal synchrony: A versatile code for the definition of relations? *Neuron*, 24, 49–65.
- Skyrms, B. (1986). *Choices and Chance*. Belmont, CA: Wadsworth.
- Sloman, A. (1999). Cognitive architecture. In R. A. Wilson and F. C. Keil (eds.), *The MIT Encyclopedia of Cognitive Science*. Cambridge, MA: MIT Press.
- Smith, L., and Thelen, E. (2003). Development as a dynamical system. *Trends in Cognitive Sciences*, 7, 343–8.
- Spelke, E. S. (1988). The origins of physical knowledge. In L. Weiskrantz (ed.), *Thought without Language*. Oxford: Oxford University Press.
- Spelke, E. S., and Kinzler, K. D. (2007). Core knowledge. *Developmental Science*, 10, 89–96.
- Spelke, E. S., and Van de Walle, G. (1993). Perceiving and reasoning about objects: Insights from infants. In N. Eilan, R. McCarthy, and B. Brewer (eds.), *Spatial Representation*. Oxford: Blackwell.
- Spelke, E. S., Gutheil, G., Van de Walle, G., Kosslyn, S. M., and Osherson, D. N. (1995). The development of object perception. In S. M. Kosslyn and D. N. Osherson (eds.), *An Invitation to Cognitive Science: Vol. 2. Visual Cognition* (2nd edn.). Cambridge, MA: MIT Press.
- Spencer, J. P., Austin, A., and Schutte, A. R. (2012). Contributions of dynamic systems theory to cognitive development. *Cognitive Development*, 27, 401–18.

- Spencer, J. P., Perone, S., and Buss, A. T. (2011). Twenty years and going strong: A dynamic systems revolution in motor and cognitive development. *Child Development Perspectives*, 5, 260–6.
- Spencer, J. P., Thomas, M. S. C., and McClelland, J. L. (2009). *Toward a Unified Theory of Development: Connectionism and Dynamic Systems Theory Reconsidered*. New York: Oxford University Press.
- Sperber, D., Cara, F., and Girotto, V. (1995). Relevance theory explains the selection task. *Cognition*, 57, 31.
- Sperling, G. (1960). The information available in brief visual presentations. *Psychological Monographs*, 74, 1–29.
- Spivey, M. (2007). *The Continuity of Mind*. New York: Oxford University Press.
- Stein, J. F., and Stoodley, C. S. (2006). *Neuroscience: An Introduction*. Oxford: Oxford University Press.
- Sterelny, K. (1990). *The Representational Theory of Mind*. Oxford: Blackwell.
- Sullivan, J. A. (2009). The multiplicity of experimental protocols: A challenge to reductionist and non-reductionist models of the unity of neuroscience. *Synthese*, 167, 511–39.
- Sun, R. (ed.) (2008). *The Cambridge Handbook of Computational Psychology*. Cambridge: Cambridge University Press.
- Sutton, R. S., and Barto, A. G. (1998). *Reinforcement Learning: An Introduction*. Cambridge MA: MIT Press.
- Tamir, D. I., and Mitchell, J. P. (2010). Neural correlates of anchoring-and-adjustment during mentalizing. *Proceedings of the National Academy of Sciences USA*, 107, 10827–32.
- Thelen, E., and Smith, L. (eds.) (1993). *A Dynamical Systems Approach to the Development of Cognition and Action*. Cambridge, MA: MIT Press.
- Thelen, E., Schöner, G., Scheier, C., and Smith, L. B. (2001). The dynamics of embodiment: A field theory of infant perseverative reaching. *Behavioral and Brain Sciences*, 24, 1–86.
- Thigpen, N. N., and Keil, A. (2017). Event-related potentials. Reference Module in Neuroscience and Biobehavioral Psychology: Elsevier. [online resource]
- Thomas, J. B., Brier, M. R., Bateman, R. J., Snyder, A. Z., Benzinger, T. L., Xiong, C., ... Ances, B. M. (2014). Functional connectivity in autosomal dominant and late-onset Alzheimer disease. *JAMA Neurology*, 71(9), 1111–22.
- Thompson, S. P., and Newport, E. L. (2007). Statistical learning of syntax: The role of transitional probability. *Language Learning and Development*, 3, 1–42.
- Tolman, E. C. (1948). Cognitive maps in rats and men. *Psychological Review*, 55, 189–208.
- Tolman, E. C., and Honzik, C. H. (1930). "Insight" in rats. *University of California Publications in Psychology*, 4, 215–32.
- Tolman, E. C., Ritchie, B. F., and Kalish, D. (1946). Studies in spatial learning, II: Place learning versus response learning. *Journal of Experimental Psychology*, 36, 221–9.
- Tononi, G., and Koch, C. (2008). The neural correlates of consciousness. *Annals of the New York Academy of Sciences*, 1124, 239–61.
- Trappenberg, T. (2010). *Fundamentals of Computational Neuroscience* (2nd edn.). Oxford: Oxford University Press.
- Trauble, B., Marinovic, V., and Pauen, S. (2010). Early theory of mind competencies: Do infants understand others' beliefs? *Infancy*, 15(4), 434–44.
- Trevethan, C. T., Sahraie, A., and Weiskrantz, L. (2007). Form discrimination in a case of blindsight. *Neuropsychologia*, 45, 2092–103.
- Tsotsos, J. K. (2011). *A Computational Perspective on Visual Attention*. Cambridge, MA: MIT Press.
- Tulving, E. (1972). Episodic and semantic memory. In E. Tulving and W. Donaldson (eds.), *Organization of Memory*. New York: Academic Press.

- Turing, A. M. (1936–7). On computable numbers: With an application to the Entscheidungsproblem [Decision Problem]. *Proceedings of the London Mathematical Society*, 42, 3–4.
- Turing, A. M. (1950). Computing machinery and intelligence. *Mind*, 59, 433–60.
- Tye, M. (1991). *The Imagery Debate*. Cambridge, MA: MIT Press.
- Umlita, M. A., Kohler, E., Gallese, V., et al. (2001). I know what you are doing: A neurophysiological study. *Neuron*, 31, 155–65.
- Ungerleider, L. G., and Mishkin, M. (1982). Two cortical visual systems. In D. J. Ingle, R. J. W. Mansfield, and M. A. Goodale (eds.), *Analysis of Visual Behavior*. Cambridge, MA: MIT Press.
- Vaina, L. M. (ed.) (1991). *From the Retina to the Neocortex*. Boston: Springer.
- Van den Bussche, E., Hughes, G., Humbeeck, N. V., and Reynvoet, B. (2010). The relation between consciousness and attention: An empirical study using the priming paradigm. *Consciousness and Cognition*, 19, 86–9.
- Van Essen, D. C., and Gallant, J. L. (1994). Neural mechanisms of form and motion processing in the primate visual system. *Neuron*, 13, 1–10. Reprinted in W. Bechtel, P. Mandik, J. Mundale, and R. S. Stufflebeam (eds.) (2001), *Philosophy and the Neurosciences: A Reader*. Malden, MA: Blackwell.
- Van Gelder, T. (1995). What might cognition be, if not computation? *The Journal of Philosophy*, 92, 345–81.
- Van Gelder, T. (1998). The dynamical hypothesis in cognitive science. *Behavioral and Brain Sciences*, 21, 615–28.
- Voyer, D., Voyer, S., and Bryden, M. P. (1995). Magnitude of sex differences in spatial abilities: A meta-analysis and consideration of critical variables. *Psychological Bulletin*, 117, 250–70.
- Wagner, A. D., Schacter, D. L., Rotte, M., Koutstaal, W., Maril, A., Dale, A. M., ... Buckner, R. L. (1998). Building memories: Remembering and forgetting of verbal experiences as predicted by brain activity. *Science*, 281(5380), 1188–91.
- Wang, S.-H., and Baillargeon, R. (2008). Detecting impossible changes in infancy: A three-system account. *Trends in Cognitive Sciences*, 12, 17–23.
- Warrington, E., and Taylor, A. M. (1973). The contribution of the right parietal lobe to object recognition. *Cortex*, 9, 152–64.
- Warrington, E., and Taylor, A. M. (1978). Two categorical stages of object recognition. *Perception*, 7, 695–705.
- Warwick, K. (2012). *Artificial Intelligence: The Basics*. London: Routledge.
- Watson, J. B. (1913). Psychology as the behaviorist sees it. *Psychological Review*, 20, 158–77.
- Waytz, A., and Mitchell, J. (2011). Two mechanisms for simulating other minds: Dissociations between mirroring and self-projection. *Current Directions in Psychological Science*, 20, 197–200.
- Webb, B. (1995). Using robots to model animals: A cricket test. *Robotics and Autonomous Systems*, 16(2), 117–34.
- Weisberg, D. S. (2015). Pretend play. *Wiley Interdisciplinary Reviews: Cognitive Science*, 6(3), 249–61.
- Weiskopf, D. A. (2004). The place of time in cognition. *British Journal for the Philosophy of Science*, 55, 87–105.
- Westermann, G., and Ruh, N. (2012). A neuroconstructivist model of past tense development and processing. *Psychological Review*, 119, 649–67.
- White R. L., III, and Snyder, L. H. (2007). Subthreshold microstimulation in frontal eye fields updates spatial memories. *Experimental Brain Research*, 181, 477–92.
- Wicker, B., Keysers, C., Plailly, J., Royet, J. P., Gallese, V., and Rizzolatti, G. (2003). Both of us disgusted in my insula: The common neural basis of seeing and feeling disgust. *Neuron*, 40, 655–64.

- Wilson, R. A. (2008). The drink you have when you're not having a drink. *Mind and Language*, 23, 273–83.
- Wimmer, H., and Perner, J. (1983). Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children's understanding of deception. *Cognition*, 13, 103–28.
- Winfield, A. F. T. (2012). *Robotics: A Very Short Introduction*. Oxford: Oxford University Press.
- Winograd, T. (1972). *Understanding Natural Language*. New York: Academic Press.
- Winograd, T. (1973). A procedural model of language understanding. In R. C. Schank and A. M. Colby (eds.), *Computer Models of Thought and Language*. San Francisco: W. H. Freeman.
- Womelsdorf, T., Schoffelen, J. M., Oostenveld, R., et al. (2007). Modulation of neuronal interactions through neuronal synchronization. *Science*, 316, 1609–12.
- Woodward, A., and Needham, A. (2009). *Learning and the Infant Mind*. Oxford: Oxford University Press.
- Wu, X., Kumar, V., Quinlan, J. R., et al. (2008). Top 10 algorithms in data mining. *Knowledge and Information Systems*, 14, 1–37.
- Xu, F., and Tenenbaum, J. B. (2007). Word learning as Bayesian inference. *Psychological Review*, 114(2), 245–72.
- Zacks, J. M. (2008). Neuroimaging studies of mental rotation: A meta-analysis and review. *Journal of Cognitive Neuroscience*, 20, 1–19.
- Zeki, S. M. (1978). Functional specialization in the visual cortex of the rhesus monkey. *Nature*, 274, 423–8.
- Zelazo, P. D., Moscovitch, M., and Thompson, E. (eds.) (2007). *The Cambridge Handbook of Consciousness*. Cambridge: Cambridge University Press.
- Zylberberg, A., Dehaene, S., Roelfsema, P. R., and Sigman, M. (2011). The human Turing machine: A neural framework for mental programs. *Trends in Cognitive Sciences*, 15, 293–300.



## INDEX FOR COGNITIVE SCIENCE (3RD EDITION)

- A-consciousness (access consciousness) 393–4,  
    444
- A-Not-B error, dynamical systems approach  
    161–5
- abduction (abductive reasoning) 444
- Abrams, Richard 392–3
- absolute judgment 444
- access consciousness (A-consciousness) 393–4,  
    444
- action potentials 9, 237, 444
- activation function 444
- activation space 301
- ACT-R cognitive architecture 220–2  
    basic structure and features 220–2  
    hybrid architecture 222–4
- ADHD 439
- Adolphs, Ralph 370
- agent, definition 204
- agent architectures 204–7  
    definition 204  
    goal-based agents 205–6  
    hybrid architectures 219–24  
    learning agents 206–7  
    massive modularity hypothesis 210–19  
    modularity of mind (Fodor) 207–10  
    reflex agents 204–5  
    simple reflex agents 204–5
- algorithms 106, 444  
    concept 22–5
- local learning algorithms 140
- SHRDLU program 40–7
- in transformational grammar 27–8
- Allen robot 424–7
- AlphaGo program 327–30
- AlphaGo Zero program 329–30
- Alzheimer’s disease 439
- Amazon’s Alexa 318
- analytic tradition 5
- anatomical connectivity 444
- AND Boolean function 129–31
- anterograde amnesia 444
- anthropology, role in cognitive science 3–5
- architectures *see* agent architectures
- Arnheim, Rudolf 180
- artificial agents  
    architectures for 204–7  
    types of 204
- artificial intelligence (AI) 47, 100  
    chatbots 39–40  
    role in cognitive science 3–5  
    SHRDLU program 40–7  
    strong and weak AI 118  
    voice-based 318  
    *see also* machine learning; robotics
- artificial neural networks (connectionist  
    networks) 77–8, 124–5, 444  
    ability to learn from “experience” 143  
    biological plausibility 139–41  
    distributed representations 141  
    Fodor–Pylyshin objection to neural network  
        models 301–2  
    information processing 78–80  
    information storage and information  
        processing not clearly distinguished 142–3  
    key features of information processing 141–3  
    models of children’s physical reasoning  
        293–300  
    models of language learning 266–74  
    multilayer networks 136–41  
    neurons and network units 125–8  
    pattern recognition 78–80  
    relationship to physical symbol models  
        300–2  
    single-layer networks 128–36
- Aslin, Richard 275
- association 17
- attention  
    locus of selection problem 241–6  
    networks of attention 246–9
- attention effect 243
- attention selection models 241–2
- attractors 444
- autism, pretend play and 341–2
- autoencoders 322–4, 444
- autonomous vehicles 442–3
- Baars, Bernard 396–7
- backpropagation algorithm 444
- Baillargeon, Renée 287–8, 347–8
- balance beam problem, network modeling  
    297–300
- Banaji, Mazharin 374
- Baron-Cohen, Simon 342–6, 349–53, 373
- basins of attraction in state space 157–8
- Bayes, Thomas 172
- Bayesian language learning 274–80

- Bayesianism 444  
basic elements 172–9  
binocular rivalry (case study) 182–6  
conditional probability 175–6  
decision-making 189–90  
degrees of belief 173  
expected utility theory 187–90  
neuroeconomics 186–98  
neurons that code for expected utility 190–8  
perception as a Bayesian problem 179–86  
probability calculus 173–4  
subjective probability 173–4  
Bayes's Rule 176–9, 444  
application to language learning 274–5  
likelihood of the evidence 178–9  
posterior probability 178–9  
prior probability 178–9  
behavior-based robotics 427–32, 444  
behavioral finance 186  
behaviorism 444  
reaction against 16–22  
belief–desire psychology 107–8  
belief system 210  
Bengio, Yoshua 318, 320–1, 325  
Berger, Theodore 440  
Bergson, Henri 401  
Bernoulli, Daniel 189  
Bernoulli, Nicholas 189  
Berti, Anna 387  
binary Boolean functions 128–31  
binding problem 444  
binocular rivalry 182–6, 445  
biorobotics 418–23, 445  
Bisiach, Eduardo 385  
bits 28, 50, 445  
blindsight 445  
nonconscious processing 384–7  
what is missing 389  
Block, Ned 393–4  
BOLD (blood oxygen level dependent) fMRI  
signal 86–7, 240, 445  
neural correlates 90–2  
relation to cognition 250–1  
Boolean functions 445  
single-layer networks 128–36  
bots 204  
bottom-up approach to study of the mind  
70–6  
brain  
activity associated with remembering visual  
experiences 87–9  
default mode network (DMN) 439  
descriptive vocabulary 72  
functions of the lobes 68–70  
brain anatomy 68–70  
Brodmann areas 69–70  
brain atlases 252  
brain connectivity, Human Connectome  
Project 438–9  
brain-damaged patients, cognitive problems 73  
brain imaging *see* neuroimaging  
brain mapping  
anatomical connectivity 232–7  
attention selection models 241–2  
blood flow and blood oxygen levels 240–1  
combining ERPs and single-unit recording  
242–6  
combining resources in networks of attention  
246–9  
combining resources in the locus of selection  
problem 241–6  
EEG (electroencephalography) 237–41  
electrical activity of the brain 237–40  
fMRI (functional magnetic resonance  
imaging) 240–1  
functional connectivity versus effective  
connectivity 252–3  
functions of the lobes 232–4  
hypotheses about visuospatial attention  
248–9  
MEG (magnetoencephalography) 240–1  
from neuroimaging data to maps 249–53  
neuroscientific techniques 237–41  
noise in the neuroimaging system 251–2  
PET (positron emission tomography) 240–1  
relation between blood flow and cognition  
250–1  
structure and function in the brain 230–7  
tract tracing 234  
Broadbent, Donald 28–32, 241–2, 397  
Brodmann, Korbinian 232–4  
Brodmann areas 70, 232–4, 445  
Brooks, Rodney 416–18, 423–8  
Bruner, Jerome 180  
Busemeyer, Jerome R. 161  
bytes 50  
CAPTCHA tool 316  
Carruthers, Peter 396  
causation by content 108–10  
central cognitive processing 209–10  
cerebral cortex 445  
Chalmers, David 394–6  
Changeux, Jean-Pierre 398–400  
channel capacity 30, 445  
chatbots 39–40, 445  
cheater detection module 211–13, 445  
Cheng, Gordon 440  
Chilausky, Richard 314–15  
child development, dynamical systems  
approach 158–67  
Chinese room argument 315, 452  
Chomsky, Noam 26–8, 38, 260, 265  
chunking information 30, 445

- Church, Alonzo 24–5  
 Church–Turing thesis 24–5, 445  
 Churchland, Patricia 5  
 classical/Pavlovian conditioning 16–17, 445  
 cochlear implants 440  
 cocktail party phenomenon 31–2  
 cognitive-deductive system 41–2  
 cognitive maps 21  
 cognitive model of single-word processing 82–4  
 cognitive neuroscience of mindreading 365–75  
 cognitive processes  
     central processing 209–10  
     massive modularity hypothesis 210–19  
     modular processing 208–9  
 cognitive psychology 47–8  
 cognitive science  
     aim of 10  
     functional view 66–8  
     interdisciplinary nature 3–5  
     range and scope 5–9  
     space of 10–11  
     three dimensions of variation 10  
     unified Theory of Cognition 10  
 cohesion, principle of 288–90, 451  
 Colby, Carol 248–9  
 competitive networks 140–1, 445  
 complex behaviors  
     planning and organization 21–2  
     problem of serial order 21–2  
 computation 445  
 computation theory 22–5  
 computational governor 154–6  
 computational model of motor control 159  
 computational modeling of the brain 76–80  
 computational neuroscience 124–5, 445  
 conditional probability 175–6, 445  
 conditioned stimulus 17  
 conditioning 16–17  
 congruence priming 445  
 connectionist approach to language learning 266–74  
 connectionist modelers 124–5  
 connectionist models of tense learning 269–74  
 connectionist networks *see* artificial neural networks  
 connectivity, anatomical 445  
 connectivity matrices 234  
 Connell, Jonathan 427  
 conscious awareness, information processing without 382–7  
 consciousness  
     access consciousness (A-consciousness) 393–4  
     blindsight and 384–7  
     challenge of understanding 380–1  
     conscious and nonconscious vision 389–92  
     diversity of research approaches 400–1  
     easy and hard problems of 394–6  
     global workspace theory of 396–400  
     inadequacy of information-processing models of the mind 380–1  
     Knowledge Argument 380–1  
     phenomenal consciousness (P-consciousness) 393–4  
     priming and 382–4  
     two types of 393–4  
     unilateral spatial neglect and 384–7  
     what is missing in blindsight and spatial neglect 389  
     what is missing in masked priming 392–3  
     what it is for 387–93  
 contact, principle of 289–90, 451  
 contention scheduling 397  
 continuation, principle of 181  
 continuity, principle of 451  
 continuity constraint 290  
 contralateral organization 445  
 convolutional neural networks (ConvNets) 324–7, 446  
 invariance under translation 326–7  
 shared weights 326  
 sparse connectivity 325–6  
 Cooper, Lynn 47–8  
 cooperation, evolution of 213–15  
 co-opted mechanisms, role in mindreading 369–75  
 co-opted systems 446  
 Corbetta, Maurizio 249  
 corpus callosum 446  
 Cosmides, Leda 211, 213–14, 216–18  
 counterfactual thinking 362, 446  
 Courville, Aaron 325  
 covert attention 246, 249, 446  
 Cox, James 212–13  
 cricket phonotaxis 419–20  
 cross-lesion disconnection experiments 73–6, 446  
 cross-talk 446  
 Damasio, Antonio 396  
 Darwinian modules 210–11  
     cheater detection module 211–13  
 Dayan, Peter 330  
 Decision Field Theory 161  
 decision-making, Bayesian approach 189–90  
 decision theory 186  
 decision trees 308–10, 446  
 declarative memory, information accessibility 223  
 deep learning 317–18, 446  
     application in autonomous vehicles 442–3  
     autoencoders 322–4  
     convolutional neural networks (ConvNets) 324–7  
     machinery of 321–7  
     visual cortex and 318–21

- deep reinforcement learning 327–30  
deep structure (phrase structure) of a sentence 26–8, 446  
default mode network (DMN) 439  
Dehaene, Stanislas 396–400  
delayed saccade tasks 248–9  
della Porta, Giambattista 182  
delta rule 450  
Dennett, Daniel 401  
deontic conditionals 211  
dichotic listening experiments 31–2, 446  
diffusion tractography 236  
digital information storage 50–1  
dimensionality reduction 322  
diminishing marginal utility 189  
dishabituation paradigm 286–92, 446  
distributed representations in neural networks 141, 446  
domain-general mental architecture, arguments against 214–18  
domain-specific mechanisms 446  
domains 128  
dopamine neurotransmitter activity 330  
dorsal visual pathway 70–6, 446  
double dissociation 385, 446  
Draine, Sean 392–3  
drawbridge experiments 287–8  
Duncan, John 242  
dynamic field model 163–5, 167  
dynamical modeling 150–3  
dynamical systems, definition 150–3  
dynamical systems approach  
    A-Not-B error 161–5  
    applications in child development 158–67  
    assessment of the approach 166–7  
    modeling motor control 159–61  
dynamical systems hypothesis 446  
dynamical systems theory 446  
    approach to cognitive science 149–58  
    basins of attraction in state space 157–8  
    relevance to cognitive science 153–8
- early selection model of attention 241–2, 446  
Ebbinghaus illusion 391  
EEG (electroencephalography) 446  
    brain mapping 237–41  
effective connectivity 446  
effector systems 204  
ELIZA program 39–40  
Elman, Jeff 266  
emergent behavior 427, 431–2  
entropy 446  
EPIC architecture 222  
epilepsy 439–40  
Evarts, Edward 190  
event-related fMRI 86–9, 447  
event-related magnetic fields 447  
event-related potentials (ERPs) 240, 447  
evolution of cooperation 213–15  
evolutionarily stable strategy 214  
exoskeletons (robot suits) 440  
expected utility 447  
    concept 189–90  
    neurons that code for 190–8  
    role of the lateral intraparietal (LIP) neurons 191–8  
    saccadic eye movement experiments 191–8  
    theory 187–90  
expected value 188–9  
expert systems  
    decision trees 308–10  
    machine learning and 308–15  
    research 447  
eyewitness testimony 441
- factive states 343  
false belief  
    implicit and explicit understanding of 347–8  
    Perner's model of theory of mind development 360–3  
selection processor hypothesis 358–60  
why it takes so long to understand 358–63  
false belief task 447  
    used to study mindreading 342–6  
false photograph task 369  
Fang, Fang 391  
feature engineering 316–17, 447  
feature learning 317 *see also representation learning*  
feedforward networks 137, 447  
Felleman, Daniel J. 234  
first-order predicate calculus 408  
fixed neural architectures 447  
flowchart model 28–9  
fMRI (functional magnetic resonance imaging) 240–1, 447  
    brain activity associated with remembering visual experiences 87–9  
    event-related 86–9  
    functional neuroimaging 84–7  
    *see also BOLD fMRI signal*  
Fodor, Jerry 106–14, 207–10, 263–5, 301–2  
folk physics 288, 443, 447  
formal property 447  
fovea 447  
frame problem 447  
Frege, Gottlob 261  
Friston, Karl 91, 182–6  
Frith, Uta 342–6  
Frost, Douglas 386  
functional connectivity 447  
functional decomposition 447

- functional neuroimaging 447  
   with fMRI 84–7  
   with PET 81–4
- functional systems 447
- functional view of cognitive systems 66–8
- functions 128–30
- future challenges and opportunities 438–43  
   autonomous vehicles 442–3  
   brain connectivity 438–9  
   default mode network (DMN) 439  
   Human Connectome Project 438–9  
   law and cognitive science 441–2  
   neural prosthetics 440  
   what the brain is doing when it appears not to be doing anything 439
- Gahl, Susan 277–8
- General Problem Solver program 100, 105–6
- Gestalt school of perceptual psychology 180–1
- Glimcher, Paul 191–8
- global neuronal workspace theory 398–400
- global workspace theory of consciousness 396–400, 447  
   building blocks 396–7  
   versions 397–400
- goal-based agents 205–6
- Goel, Vinod 366
- GOFAI (good old-fashioned artificial intelligence) robotics 448  
   SHAKY robot 408–14, 416
- Goldberg, Michael 191
- Goldman, Alan 363–4
- Goodale, Melvyn, theory of vision 389–92
- Goodfellow, Ian 325
- Google  
   Deep Mind research program 317, 327–30  
   self-driving cars 442
- GoogLeNet 325
- Gopnik, Alison 292–3
- Gordon, Robert 365
- Gorman, Paul 78–80
- graceful degradation 448
- Greenwald, Anthony 392–3
- Griggs, Richard 212–13
- Haldane, John Scott 401
- Halligan, Peter 387
- halting problem 23–5, 448
- Hamilton, W. D. 217–18
- hard problem of consciousness 448
- Harris, Paul 363
- Haugeland, John 415
- He, Shen 391
- Heal, Jane 365
- Hebb, Donald 131–2
- Hebbian learning 131–2, 448
- Held, Richard 386
- hemiagnosia *see* unilateral spatial neglect
- hemineglect *see* unilateral spatial neglect
- Herrnstein, Richard 196
- heuristic strategies 105, 214, 448
- hidden layers 137, 448
- hidden units 77–8, 136, 448
- Hinton, Geoffrey 318, 320–1, 324
- hippocampal prosthetics 440
- hippocampus 7
- historical landmarks in cognitive science  
   approaches to understanding information 32–4  
   computation theory and the algorithm concept 22–5  
   information-processing models in psychology 28–32  
   interdisciplinary model of vision 53–61  
   language-processing systems and micro-worlds 38–47  
   linguistics and formal analysis of language 25–8  
   reaction against behaviorism 16–22  
   representation of mental images 47–53  
   turn to the brain 66–92
- Hohwy, Jakob 182–6
- Honzik, C. H. 17–20
- Human Connectome Project 438–9
- Human Genome Project 438
- hybrid architectures 219–24
- ID3 algorithm for machine learning 310–15
- ImageNet Large-Scale Visual Recognition Challenge 324–5
- imagery debate 47–53
- infant cognition  
   connectionist models 293–300  
   dishabituation paradigm and 286–92  
   interpretation of dishabituation experiments 292–3  
   modeling object permanence 295–7  
   modeling the balance beam problem 297–300  
   neural network models 293–300  
   traditional views 286
- infant folk physics, underlying information processing 292–3
- information, concepts of 32
- information channel 448
- information processing  
   approaches to understanding 32–4  
   artificial neural networks 78–80, 141–3  
   bottleneck 30  
   channel capacity 30  
   without conscious awareness 382–7  
   formal analysis of language 25–8  
   hierarchically organized behavior 21–2  
   how it works 33  
   human limitations 29–30

- learning without reinforcement 17–20  
mental imagery 50–3  
models in psychology 28–32  
neurally inspired models 124–8  
organization of the mind and 223–4  
reaction against behaviorism 16–22  
sensory information 30–2  
spatial learning studies 20–1  
specialized systems for 33–4  
information-processing systems, levels of explanation 53–5  
information theory 28–9  
informational encapsulation 448  
innatism about language 262, 265  
insects, robotic studies 418–23  
integration challenge xxiv  
integration principle 237, 448  
intelligence, physical symbol system hypothesis 100–6  
intelligent action, and the physical symbol system 106  
intentional realism 108–10, 448  
intentionality 448  
interdisciplinary nature of cognitive science 3–5  
interocular suppression 448  
iSee 443  
isotropic property of central processing 210
- Jackson, Frank 380–1  
James, William 286  
Jasper, Herbert 190  
Jenkins, E. 298–9  
joint visual attention 351–2, 448
- Kanwisher, Nancy 368–9  
Kelly, William 374  
Kerszberg, Michel 398–400  
Keyser, Samuel Jay 5  
Kieras, David 222  
kin selection model 217–18  
Knowledge Argument 380–1, 448  
Koch, Christoph 91  
Koffka, Kurt 180  
Köhler, Wolfgang 180  
Kosslyn, Stephen 52–3  
Kuczaj, Stan 267
- language  
natural language-processing systems 38–47  
study of 25–8
- language learning  
Bayesian language learning 274–80  
challenge of tense learning 267–9  
connectionist approach 266–74  
connectionist models of tense learning 269–74  
language of thought and 263–5
- learning linguistic categories 278–80  
neural network models 266–74  
probabilities in word and phrase segmentation 275–6  
rules and language 260–2  
understanding a language and learning a language 261–2  
understanding pronouns 276–8
- language of thought hypothesis 106–14, 448  
language learning and 263–5  
main claims 110  
relation between syntax and semantics 110–14  
structure of the argument for 113–14
- language-processing systems  
ELIZA program 39–40  
SHRDLU program 40–7
- Lashley, Karl 21–2, 409
- late selection models of attention 242, 448
- latency 243
- latent learning 17–20
- lateral geniculate nucleus (LGN) 319
- lateral intraparietal area (LIP)  
delayed saccade experiments 248–9  
role of neurons in expected utility 192–8
- law, connections with cognitive science 441–2
- learning  
without reinforcement 17–20  
in single-layer networks 131–4
- learning agents 206–7
- learning algorithms 77–8, 140
- LeCun, Yann 318, 320–1
- Leibniz's Mill 448
- Leslie, Alan 342–6, 348  
model of infant pretend play 336–7  
on pretend play and metarepresentation 337–40  
selection processor hypothesis 358–60
- levels of explanation 53–5  
neuroscience and psychology compared 5–9
- lexical access 82, 449
- lexical decision task 383
- lexical processing, mapping the stages of 80–4
- Li, Fei-Fei 324
- Li, Kai 324
- Lidz, Jeffrey 277
- likelihood of the evidence 178–9, 449
- linearly separable functions 134–6, 449
- linguistics 25–8  
role in cognitive science 3–5
- local algorithm 449
- local field potential (LFP) 92, 449
- local learning algorithms 140
- localist networks 141
- locus of selection problem 241–6, 449
- logic programming 408, 413–14
- Logic Theory Machine 100

- logical consequence 449  
 logical deducibility 449  
 Logothetis, Nikos 91–2  
 Luria, Alexander 230–2  
 Luzzatti, Claudio 385
- machine learning 449  
 algorithms 449  
 AlphaGo program 327–30  
 decision trees 308–10  
 deep learning 317–18  
 deep learning and the visual cortex 318–21  
 deep reinforcement learning 327–30  
 expert systems and 308–15  
 ID3 algorithm 310–15  
 machinery of deep learning 321–7  
 representation learning 315–17  
 machine table 24  
 Macrae, Neil 374  
 macroeconomics 186  
 magnetoencephalography *see* MEG  
 mandatory application 449  
 mapping functions 128–30  
 Marchman, Victoria 272–3  
 Marcus, Gary 273  
 Marr, David 53–61, 300–1, 418  
 Marshall, John 387  
 masked priming 392–3, 449  
 massive modularity hypothesis 210–19, 449  
 argument from error 216  
 argument from statistics and learning 216–18  
 arguments against domain-general mental architecture 214–18  
 cheater detection module 211–13  
 evaluating the arguments for 218–19  
 evolution of cooperation 213–15  
 evolution of mental architecture 216–18  
 kin selection model 217–18  
 Matarić, Maja 428–32  
 matching behavior 196  
 mathematics, computation theory 22–5  
 McClelland, Jay 269–72, 293, 298–9  
 McCulloch, Warren 76, 128, 130  
 means-end analysis 105–6  
 MEG (magnetoencephalography) 449  
 brain mapping 240–1  
 melioration theory 196–8  
 Meltzoff, Andrew 292–3  
 memory  
   brain activity associated with remembering visual experiences 87–9  
   visual event studies 84–7  
 mental architecture 449  
 mental imagery  
   how images are represented 47–53  
   information processing 50–3  
 mental rotation of three-dimensional objects 48–50  
 metarepresentation 341–8, 449  
 metarepresentation, link to mindreading 341  
 metarepresentation, pretend play and 337–40  
 metarepresentation, why it takes so long to understand false belief 358–63  
 Metzler, Jacqueline 47–50  
 Meyer, David 222  
 Michalski, Ryszard 314–15  
 microcircuits 7–8  
 microeconomics 186  
 microelectrode recording of single neurons 90–2  
 micro-worlds 40–7, 449  
 Miller, George 5, 28–30  
*Miller v. Alabama* (2012) 442  
 Milner, David, theory of vision 389–92  
 Milward, T. T. 140–1  
 mindreading 443  
   autism and 341–2  
   cognitive neuroscience of 365  
   implicit and explicit understanding of false belief 347–8  
   Leslie's model of pretend play 336–40  
   link to pretend play 341  
   neuroscientific studies 365–75  
   pretend play and metarepresentation 336–41  
   role of co-opted mechanisms 369–75  
   role of simulation in high-level mindreading 373–5  
   role of simulation in low-level mindreading 369–73  
   significance of pretend play 336–7  
   using the false belief task to study 342–6  
   view of simulation theory 363–5  
   why it takes so long to understand false belief 358–63  
 mindreading system 348–53  
   first steps in mindreading 349–51  
   from dyadic to triadic interactions 351–2  
   joint visual attention 351–2  
   theory of mind mechanism (TOMM) 352–69  
 Minsky, Marvin 136  
 mirror neurons 237, 371–3, 449  
 Mishkin, Mortimer 70–6  
 Mitchell, Jason 374  
 MNI brain atlas 252  
 modular cognitive processing 208–9  
 modularity of mind (Fodor) 207–10  
 modules 449  
 Montague, Read 330  
*Montgomery v. Louisiana* (2016) 442  
 morphological computation 420–3, 449

- motor control  
computational model 159  
dynamical systems approach 159–61  
robot hand (Yokoi) 421–3
- Mountcastle, Vernon 190–1
- multiagent programming, Nerd Herd robots 430–2
- multilayer networks 136–41, 450  
backpropagation algorithm 138–9  
feedforward networks 137  
hidden layers 137  
modifying the weights 138–9  
spread of activation 137–8
- multiply realizable systems 67, 450
- Munakata, Yuko 293–7
- Musk, Elon 442
- MYCIN expert systems program 308
- Naccache, Lionel 396–7
- naïve physics 288
- nativism about language 262, 265
- natural language-processing systems 38–47
- neglect phenomenon 385 *see also blindsight; unilateral spatial neglect*
- Nerd Herd robots 430–2
- networks of attention 246–9
- neural correlates of the BOLD fMRI signal 90–2
- neural networks *see artificial neural networks (connectionist networks)*
- neural prosthetics 440
- neurally inspired models of information processing 124–8
- neuroeconomics 450  
Bayes in the brain 186–98
- neuroimaging  
functional connectivity versus effective connectivity 252–3  
investigating the theory of mind system 366–9  
techniques for studying cognitive functioning 237–41
- neurolaw 441–2
- neurological model of single-word processing 82–4
- neurons 7–9  
activation functions 126–7  
microelectrode recordings 90–2  
structure and function 125–8  
threshold for firing 125–6
- neurons that code for expected utility 190–8  
combining probability and utility 196–8  
probability-detecting neurons 193–4  
utility-detecting neurons 194–5
- neurophilosophy 5
- neuroscience  
branches of 7–9  
compared with psychology 5–9
- levels of organization 7–9  
role in cognitive science 3–5  
tools and techniques 9
- neurotransmitters 8–9, 450
- New Look perceptual psychology 180
- Newell, Allen 100–6
- Newport, Elissa 275–6
- Nicolelis, Miguel 440
- Nissl, Franz 232
- Nissl stain 232
- nonconscious information processing 382–7  
blindsight and unilateral spatial neglect 384–7
- Norman, Donald 397
- NOT Boolean function 130–1
- numerals 109
- object permanence 450  
development of 161–5  
modeling approaches 293–5  
network modeling 295–7  
and physical reasoning in infancy 286–93
- object substitution 336
- Onishi, Kristine 347–8
- operant conditioning 450
- OR Boolean function 129–31
- organism–environment complex 156
- overregularization errors 268, 450
- P-consciousness (phenomenal consciousness) 393–4, 450
- paired deficits 370
- paired-image subtraction paradigm 450
- Papert, Seymour 135
- parallel distributed processing 77–8
- parallel processing 450
- paralysis, development of exoskeletons (robot suits) 440
- Parkinson’s disease 439
- partial volume effects 251
- pattern recognition in artificial neural networks 78–80
- Pavlovian conditioning *see classical/Pavlovian conditioning*
- perception  
as a Bayesian problem 179–86  
binocular rivalry (Bayesianism case study) 182–6  
predictive challenge 179–81
- perceptron 450
- perceptron convergence rule 131–4, 450  
limits of 134–6
- Perner, Joseph 342, 360–3
- PET (positron emission tomography) 240–1, 450
- cortical anatomy of single-word processing 81–4

- phenomenal consciousness (P-consciousness) 393–4, 450
- philosophy, role in cognitive science 3–5
- phrase segmentation 276
- phrase structure (deep structure) of a sentence 26–8
- phrase structure grammar 450
- physical symbol models, relationship to neural network models 300–2
- physical symbol system 450
- physical symbol system hypothesis 100–6, 142, 450
- intelligent action and the physical symbol system 106
- language learning and 263
- language of thought 106–14
- Russian Room argument 114–19
- symbols and symbol systems 101–2
- transforming symbol structures 102–6
- physicalism (materialism) 381
- Piaget, Jean 161–2, 286
- Pinker, Steven 268–9, 272
- Pitts, Walter 76, 128, 130
- place learning 19–21
- PLANEX software 413–14
- Platt, Michael 191–8
- Plunkett, Kim 272–3
- Polizzi, Pamela 359
- Pöppel, Ernst 386
- Population-Average, Landmark and Surface-Based (PALS) brain atlas 252
- position-invariant object recognition 140–1
- posterior probability 178–9, 183, 450
- poverty of the stimulus arguments 265, 450
- pragmatics of conversation 46–7, 450
- precentral gyrus 234
- predicate calculus 111–14, 451
- prestriate cortex 74, 451
- pretend play
- autism and 341–2
  - Leslie's model 336–40
  - link to mindreading 341
  - metarepresentation and 336–41
  - significance of 336–7
- primary motor cortex 234
- primary somatosensory cortex 234
- primary visual cortex 7, 69–70, 234, 319, 451
- primary visual pathway 68–70
- priming 451
- consciousness and 382–4
- Prince, Alan 268–9, 272
- prior probability 178–9, 451
- prisoner's dilemma game 214–15, 451
- probability calculus 173–4
- probability-detecting neurons 193–4, 196
- probability theory 187
- procedural memory 223
- procedures 42–6
- production rules 205, 222–4
- pronomial anaphora 276–8
- propositional attitudes 107–8, 360, 451
- propositional logic (propositional calculus) 101–2, 451
- psychology
- compared with neuroscience 5–9
  - how it is organized 6–7
  - information-processing models 28–32
  - reaction against behaviorism 16–22
  - of reasoning 186
  - role in cognitive science 3–5
  - subfields and specializations 6–7
- psychophysics 31, 451
- Pylyshyn, Zenon 301–2
- Quine, Willard von Orman 210
- Quinean property of central processing 210
- Quinlan, Ross 310
- Raichle, Marcus 439
- range 128
- rats
- learning in maze running 17–20
  - spatial learning studies 20–1
- recurrent networks 296–7, 451
- recursive definition 451
- reduction 451
- Rees, Geraint 91
- reflex agents 204–5
- Regier, Terry 277–8
- reinforcement, learning without 17–20
- reinforcement learning 451
- representation 21, 32–3, 107–8, 451
- mental images 47–53
- representation learning 315–17, 451
- representational mind 362
- response learning 19–21
- retrograde amnesia 451
- reward prediction error hypothesis 330
- Rizzolatti, Giacomo 237, 371, 387
- robot reply (to the Chinese room argument) 452
- robotics
- Allen robot 424–7
  - alternatives to GOFAI robotics 414–23
  - behavior-based robotics 427–32
  - biorobotics 418–23
  - challenge of building a situated agent 415–16
  - emergent behaviors 427, 431–2
  - GOFAI robotics 408–14, 416
  - hybrid architectures 427
  - insect studies 418–23
  - morphological computation 420–3
  - multiagent programming 430–2
  - Nerd Herd robots 430–2

- preconditions for intelligent problem solving 416–18  
robot cricket 419–20  
robot fish called WANDA 420–1  
robot hand (Yokoi) 421–3  
robots 204  
SHAKY robot 408–14, 416  
situuated cognition and knowledge representation 416–18  
situuated cognition approach 415  
subsumption architectures 423–7  
TOTO robot 428–30  
Rock, Irving 180–1  
Roepstorff, Andreas 182–6  
Rolls, Edmund 140–1  
Rosenblatt, Frank 132–4  
Rumelhart, David 269–72, 293  
Russian room argument (Searle)  
physical symbol system hypothesis and 114–19  
responses to 117–19  
Turing Test and 117
- saccadic eye movements 246, 452  
experiments using 191–8  
Saffran, Jenny 275  
Sahraie, Arah 387  
Saxe, Rebecca 368–9  
scotoma 385  
search-spaces 103–6  
Searle, John 114–19, 315  
segregation principle 232, 452  
Sejnowski, Terrence 78–80, 330  
selection processor hypothesis 358–60, 452  
selective attention 31–2, 452  
selectivity/invariance problem 321  
self-driving cars 442–3  
self-reflection 374  
semantic priming 383–4, 452  
semantic property 452  
semantic system 41–2, 44–6  
sensory information processing 30–2  
sensory systems 204  
sentential logic 101–2  
SHAKY robot 408–14  
challenge of building a situated agent 416  
software for logic programming in STRIPS and PLANEX 413–14  
software for low-level activities and intermediate-level actions 409–13  
Shallice, Tim 397  
Shannon, Claude E. 28  
shared attention mechanism (SAM) 352  
shared weights 452  
Shepard, Roger 47–50  
Shepherd, Gordon 7–9  
shopping bots 204  
SHRDLU program 40–7, 415–16  
Siegler, Bob 298  
Simon, Herbert 100–6  
simple reflex agents 204–5  
simulation  
role in high-level mindreading 373–5  
role in low-level mindreading 369–73  
simulation theory  
radical simulationism 365, 452  
standard simulationism 363–4, 452  
view of mindreading 363–5  
single-layer networks  
Boolean functions 128–36  
learning in 131–4  
limits of perceptron convergence 134–6  
linearly separable functions 134–6  
perceptron convergence rule 131–4  
single-word processing, PET studies of cortical anatomy 81–4  
situuated cognition 415, 452  
knowledge representation 416–18  
Skinner, B. F. 17  
Skinner box 17  
Sloan Foundation's report (1978) 3–5  
Smith, Linda 159–61, 163–5  
solidity, principle of 451  
solidity constraint 290  
sonar target detection, artificial neural network 78–80  
space of cognitive science 10–11  
sparse connectivity 452  
spatial learning studies 20–1  
spatial neglect  
nonconscious processing 384–7  
what is missing 389  
spatial resolution 452  
spatial working memory 248–9  
spatially selective attention 246 *see also* visuospatial attention  
Spelke, Elizabeth 288–93  
SSS hybrid architecture 427  
St. Petersburg game 188–9  
state space 150–3, 452  
basins of attraction in 157–8  
stereotactic maps 251  
stimulus-onset asynchrony (SOA) 393  
stochastic gradient descent 329  
striate cortex 69–70 *see also* primary visual cortex  
STRIPS planner 413–14  
strong AI 118  
subconscious information processing hypothesis 22  
subcortex 452  
subsumption architectures 423–7, 452  
subsymbolic representation 223–4  
superior colliculus 191

- supervised learning 132, 452  
 SuperVision program 324–5  
 surface structure of a sentence 26, 452  
 Sylvian sulcus 230  
 symbol-grounding problem 453  
 symbol structures, transformation 102–6  
 symbols and symbol systems 101–2  
 synapses 8–9, 237, 453  
 syntactic structures 26–8  
 syntactic system 41–3  
 syntax 25  
 systems neuroscience 453  
 systems reply (to the Chinese room argument) 453
- Talairach-Tournoux brain atlas 252  
 task analysis hypothesis 22  
 temporal resolution 453  
 Tenenbaum, Joshua 278–80  
 Tesla 442  
 TESS (the empathizing system) 352–3  
 Thelen, Esther 159–61, 163–5  
 theory of mind mechanism (TOMM) 352–3, 453  
     development of 358–63  
     neuroimaging investigations 366–9  
 Thompson, Susan 276  
 threshold 453  
 TIT FOR TAT strategy 214, 453  
 Tolman, Edward 17–20  
 Tooby, John 211, 213–14, 216–18  
 top-down analysis, visual system 54–61  
 TOTO robot 428–30  
 Townsend, James T. 161  
 tract tracing 234  
 transformational grammar 26–8, 453  
 transitional probabilities 275  
 Trevethan, Ceri 387  
 truth conditions 263–5, 453  
 truth rules 263–5, 453  
 truth table 129–30  
 Turing, Alan 22–5, 106  
 Turing machines 24–5, 101–2, 106, 453  
 Turing Test, Russian Room argument 117  
 two visual systems hypothesis 70–6
- Umlita, Alessandra 373  
 unconditioned stimulus 17  
 Ungerleider, Leslie 70–6  
 unified Theory of Cognition 10  
 unilateral spatial neglect 76, 453  
     nonconscious processing 384–7
- universal Turing machine 106  
 unsupervised learning 140, 453  
 Urbach-Wiethe disease 370  
 utility concept 189–90, 453  
 utility-detecting neurons 194–6
- Van Essen, David 234, 252  
 Van Gelder, Tim 153–6, 166  
 ventral visual pathway 70–6, 453  
 VisNet model 140–1  
 visual cortex, deep learning and 318–21 *see also primary visual cortex*  
 visual experience, brain activity that predicts remembering 87–9  
 visual perception, binocular rivalry (Bayesianism case study) 182–6  
 visual system  
     dorsal and ventral pathways 70–6  
     interdisciplinary model 53–61  
     levels of explanation 53–5  
     Marr's model 53–61  
     position-invariant object recognition 140–1  
     top-down analysis 54–61  
     two visual systems hypothesis 70–6  
     vision for action and vision for perception 389–92  
 visuospatial attention  
     hypotheses about how it works 248–9  
     networks of attention 246–9  
 von Helmholtz, Hermann 180
- Walker, Edward 5  
 WANDA (robot fish) 420–1  
 Wason selection task 211–13, 453  
 Watt, James 153–6  
 Watt governor 153–6  
 weak AI 118  
 Webb, Barbara 419–20  
 Weiskrantz, Larry 387  
 well-formed formula 453  
 Wertheimer, Max 180  
 Wimmer, Heinz 342  
 Winograd, Terry 40–7, 415–16  
 Wittgenstein, Ludwig 262  
 word segmentation 275
- Xu, Fei 278–80
- Yokoi, Hiroshi 421–3
- Zaitchik, Debbie 369