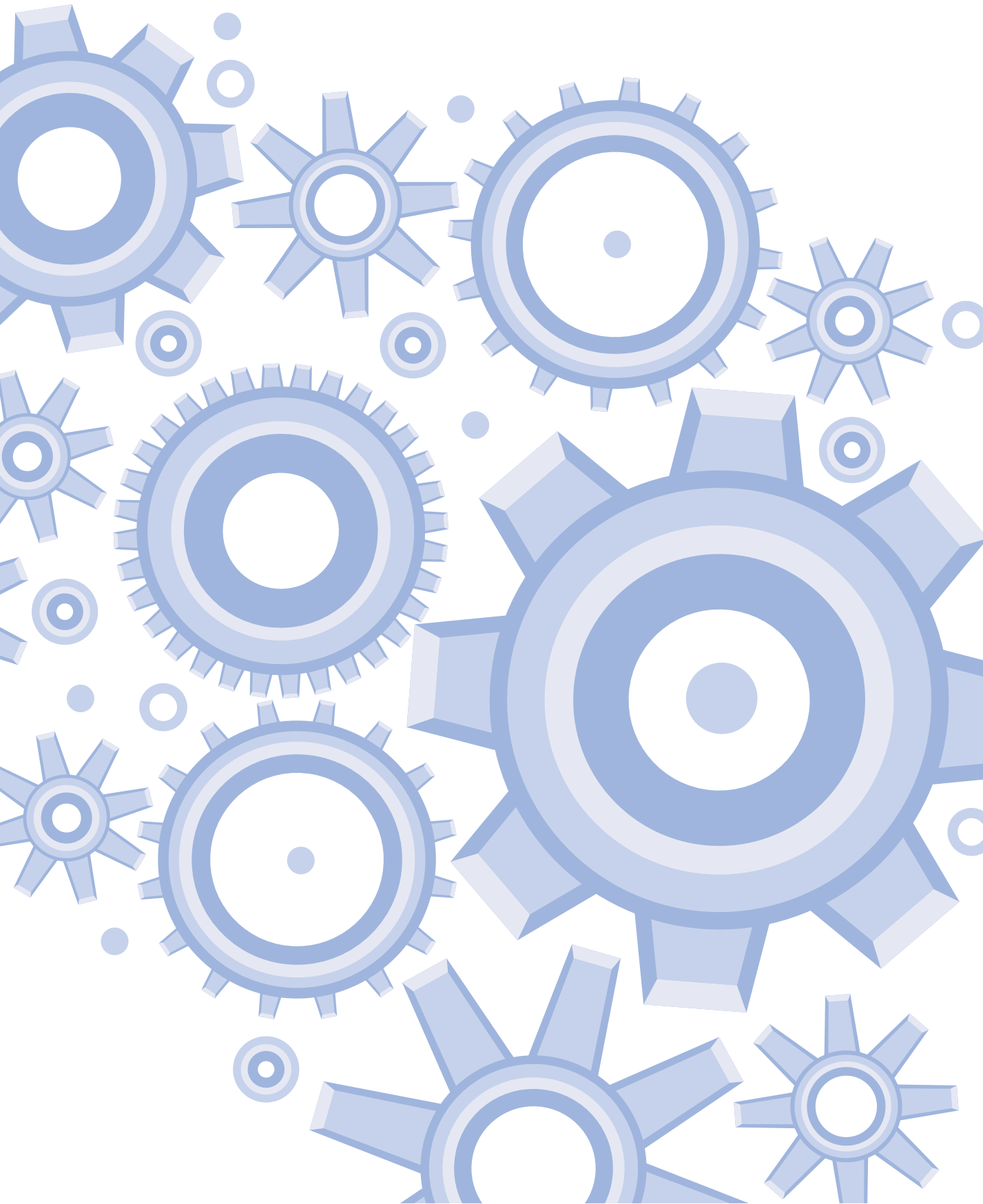# COGNITIVE SCIENCE
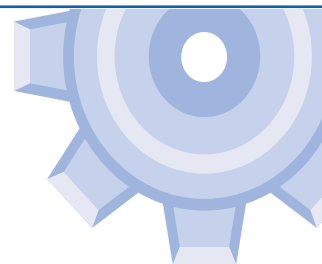
## Second edition

*Cognitive Science* combines the interdisciplinary streams of cognitive science into a unified narrative in an all-encompassing introduction to the field. This text presents cognitive science as a discipline in its own right, and teaches students to apply the techniques and theories of the cognitive scientist's "toolkit" – the vast range of methods and tools that cognitive scientists use to study the mind. Thematically organized, rather than by separate disciplines, *Cognitive Science* underscores the problems and solutions of cognitive science, rather than those of the subjects that contribute to it – psychology, neuroscience, linguistics, etc. 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 updated and enhanced to incorporate new studies and key experiments since the first edition. A new chapter on consciousness has been added.

JOSÉ LUIS BERMÚDEZ is Dean of the College of Liberal Arts and Professor of Philosophy at Texas A&M University. He has been involved in teaching and research in cognitive science for over twenty years, and is very much involved in bringing an interdisciplinary focus to cognitive science through involvement with conference organization and journals. His 100+ publications include the textbook *Philosophy of Psychology: A Contemporary Introduction* (2005) and a companion collection of readings, *Philosophy of Psychology: Contemporary Readings* (2007). He has authored the monographs *The Paradox of Self-Consciousness* (1998), *Thinking without Words* (2003), and *Decision Theory and Rationality* (2009) in addition to editing a number of collections including *The Body and the Self* (1995), *Reason and Nature* (2002), and *Thought, Reference, and Experience* (2005).
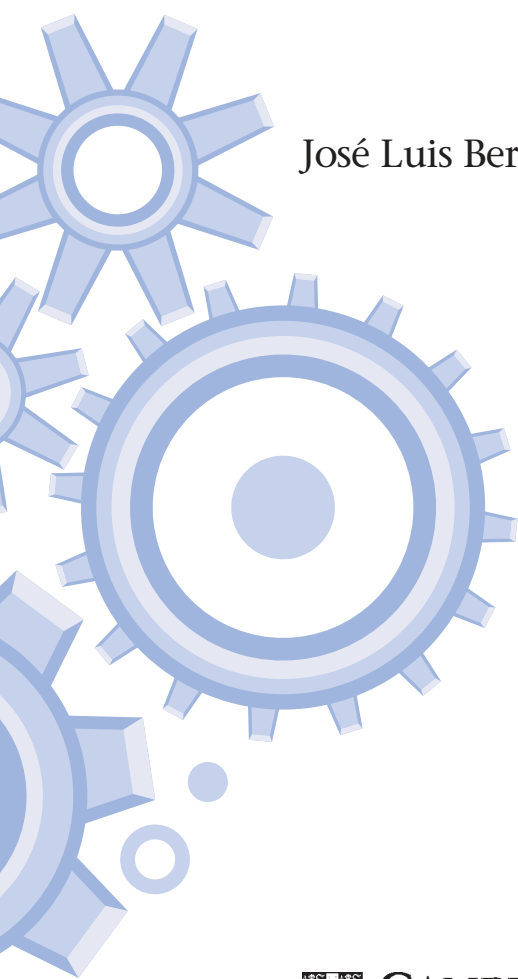
# COGNITIVE SCIENCE

## An Introduction to the Science of the Mind

### Second Edition

José Luis Bermúdez

CAMBRIDGE
UNIVERSITY PRESS

## CAMBRIDGE
### UNIVERSITY PRESS

# CONTENTS

# CONTENTS

# BOXES

# FIGURES

# TABLES

# PREFACE

## About this book

There are few things more fascinating to study 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, there are very few academic areas that 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.

This enormous range of information out there about the mind can be overwhelming, both for students and for instructors. I had direct experience of how challenging this can be when I was Director of the Philosophy–Neuroscience–Psychology program at Washington University in St. Louis. My challenge was to give students a broad enough base while at the same time bringing home that cognitive science is a field in its own right, separate and distinct from the disciplines on which it draws. I set out to write this book because my colleagues and I were unable to find a book that really succeeds in doing this.

Different textbooks have approached this challenge in different ways. Some 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 learnt much from these books and they can be very helpful. But I often have the feeling that students 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

**xxvii**

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, in order to give students a sense of cognitive science's breadth and range.

Cognitive science has only been with us for forty or so years. But in that time it has changed a lot. At one time cognitive science was associated with the idea that we can understand the mind without worrying about its biological machinery – we can understand the software without understanding the hardware, to use a popular image. But this is now really a minority view. Neuroscience is now an absolutely fundamental part of cognitive science. Unfortunately this has not really been reflected in textbooks on cognitive science. This book presents a more accurate picture of how central neuroscience is to cognitive science.

## How the book is organized

This book is organized into five parts.

### Part I: Historical overview

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: The integration challenge

The two chapters in Part II bring out what is distinctive about cognitive science. They do this in terms of what I call the integration challenge. This is the challenge of developing a unified framework that makes explicit the relations between the different disciplines on which cognitive science draws and the different levels of organization that it studies. In Chapter 4 we look at two examples of *local integration*. The first example explores how evolutionary psychology has been used to explain puzzling data from human decision-making, while the second focuses on what exactly it is that is being studied by techniques of neuro-imaging such as functional magnetic resonance imaging (fMRI).

In Chapter 5 I propose that one way of answering the integration challenge is through developing models of mental architecture. A model of mental architecture includes

1  an account of how the mind is organized into different cognitive systems, and
2  an account of how information is processed in individual cognitive systems.

This approach to mental architecture sets the agenda for the rest of the book.

## Part III: Information-processing models of the mind

The four chapters in Part III explore the two dominant models of information processing in contemporary cognitive science. The first model 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. The theoretical case for the physical symbol system hypothesis is discussed in Chapter 6, while Chapter 7 gives three very different examples of research within that paradigm – from data mining, artificial vision, and robotics.

The second model of information processing derives from models of artificial neurons in computational neuroscience and connectionist artificial intelligence. Chapter 8 explores the motivation for this approach and introduces some of the key concepts, while Chapter 9 shows how it can be used to model aspects of language learning and object perception.

## Part IV: How is the mind organized?

A mental architecture includes a model both of information processing and of how the mind is organized. The three chapters in Part IV look at different ways of tackling this second problem. Chapter 10 examines the idea that some forms of information processing are carried out by dedicated cognitive modules. It looks also at the radical claim, proposed by evolutionary psychologists, that the mind is simply a collection of specialized modules. In Chapter 11 we look at how techniques such as functional neuroimaging can be used to study the organization of the mind. Chapter 12 shows how the theoretical and methodological issues come together by working 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).

## Part V: New horizons

As emerges very clearly in the first four parts of the book, cognitive science is built around some very basic theoretical assumptions – and in particular around the assumption that the mind is an information-processing system. In Chapter 13 we look at two ways in which cognitive scientists have proposed extending and moving beyond this basic assumption. One of these research programs is associated with the dynamical systems hypothesis in cognitive science. The second is opened up by the situated/embodied cognition movement. Chapter 14 explores 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.
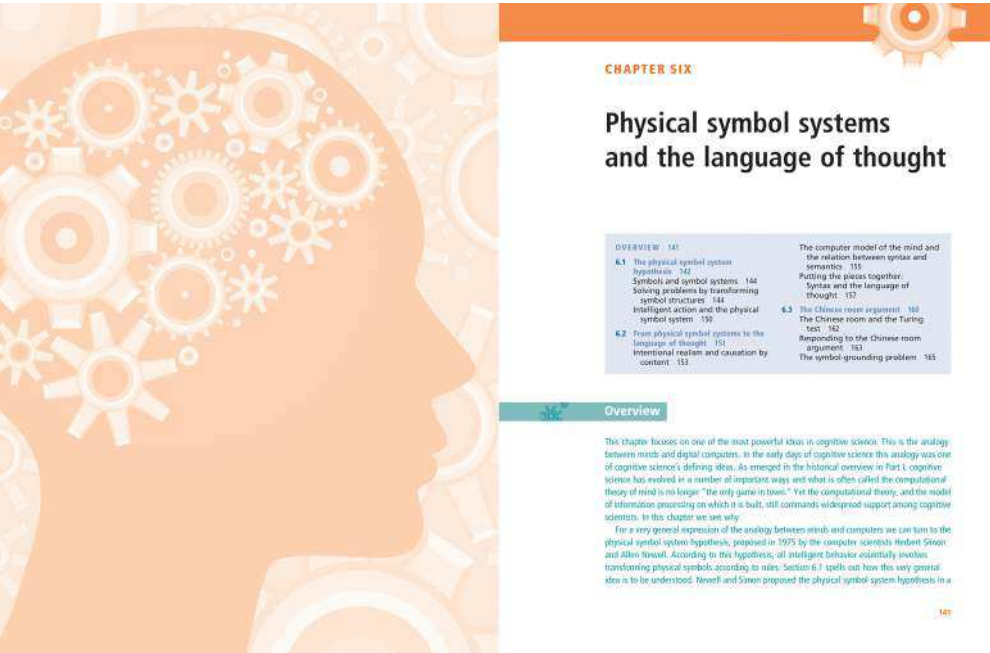
## Using this book in courses

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:

■ **Part-openers and chapter overviews**   The book is divided into five parts, as described above. Each part begins with a short introduction to give the reader a broad picture of what lies ahead. Each chapter begins with an overview to orient the reader.



■ **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).



■ **Boxes and optional material**   Boxes have been included to provide further information about the theories and research discussed in the text. Some of the more technical material has been placed in boxes that are marked optional. Readers are encouraged to work through these, but the material is not essential to the flow of the text.