# Reasoning about Knowledge

# Reasoning about Knowledge

Ronald Fagin
Joseph Y. Halpern
Yoram Moses
Moshe Y. Vardi

*To Susan, who is as happy and amazed as I am that The Book is finally completed; to Josh, Tim, and Teddy, who are impressed that their father is an Author; and to the memories of my mother Maxine, who gave me a love of learning, and of my father George, who would have been proud.*

R. F.

*To Gale, for putting up with this over the years; to David and Sara, for sometimes letting Daddy do his work; and to my mother Eva, to whom I can finally say "It's done!"*

J. Y. H.

*To my father Shimon, to Yael, Lilach and Eyal, and to the memory of my mother Ala and my brother Amir. With Love.*

Y. M.

*To Pam, who listened for years to my promises that the book is 90% done; to Aaron, who, I hope, will read this book; to my parents, Ziporah and Pinkhas, who taught me to think; and to my grandparents, who perished in the Holocaust.*

הביטו וראו, אם־יש מכאוב כמכאובי.

*"Behold and see, if there be any sorrow like unto my sorrow."*

M. Y. V.

# Contents

# Preface to the Hardcover Edition

As its title suggests, this book investigates reasoning about knowledge, in particular, reasoning about the knowledge of agents who reason about the world and each other's knowledge. This is the type of reasoning one often sees in puzzles or Sherlock Holmes mysteries, where we might have reasoning such as this:

> If Alice knew that Bob knew that Charlie was wearing a red shirt, then Alice would have known that Bob would have known that Charlie couldn't have been in the pantry at midnight. But Alice didn't know this . . .

As we shall see, this type of reasoning is also important in a surprising number of other contexts. Researchers in a wide variety of disciplines, from philosophy to economics to cryptography, have all found that issues involving agents reasoning about other agents' knowledge are of great relevance to them. We attempt to provide here a framework for understanding and analyzing reasoning about knowledge that is intuitive, mathematically well founded, useful in practice, and widely applicable.

The book is almost completely self-contained. We do expect the reader to be familiar with propositional logic; a nodding acquaintance with distributed systems may be helpful to appreciate some of our examples, but it is not essential. Our hope is that the book will be accessible to readers from a number of different disciplines, including computer science, artificial intelligence, philosophy, and game theory. While proofs of important theorems are included, the non-mathematically-oriented reader should be able to skip them, while still following the main thrust of the book.

We have tried to make the book modular, so that, whenever possible, separate chapters can be read independently. At the end of Chapter 1 there is a brief overview of the book and a table of dependencies. Much of this material was taught a number of times by the second author in one-quarter courses at Stanford University and

by the third author in one-semester courses at the Weizmann Institute of Science. Suggestions for subsets of material that can be covered can also be found at the end of Chapter 1.

Many of the details that are not covered in the main part of the text of each chapter are relegated to the exercises. As well, the exercises cover material somewhat tangential—but still of interest!—to the main thrust of the chapter. We recommend that the reader at least look over all the exercises in each chapter. Far better, of course, would be to do them all (or at least a reasonable subset). Problems that are somewhat more difficult are marked with ∗, and even more difficult problems are marked with ∗∗.

Each chapter ends with a section of notes. These notes provide references to the material covered in each chapter (as well as the theorems that are stated but not proved) and, occasionally, more details on some points not covered in the chapter. The references appearing in the notes are to the latest version of the material we could find. In many cases, earlier versions appeared in conference proceedings. The dates of the references that appear in the notes therefore do not provide a chronological account of the contributions to the field. While we attempt to provide reasonably extensive coverage of the literature in these notes, the field is too large for our coverage to be complete. We apologize for the inadvertent omission of relevant references.

The book concludes with a bibliography, a symbol index, and an index.

Many people helped us in many ways in the preparation of this book, and we are thankful to all of them. Daphne Koller deserves a very special note of thanks. She did a superb job of proofreading the almost-final draft of the book. Besides catching many typographical errors, she gave us numerous suggestions on improving the presentation in every chapter. We are very grateful to her. We would also like to thank Johan van Benthem, Adam Grove, Vassos Hadzilacos, Lane Hemaspaandra and the students of CS 487 at the University of Rochester, Wil Janssen, Hector Levesque, Murray Mazer, Ron van der Meyden, Jan Pachl, Karen Rudie, Ambuj Singh, Elias Thijsse, Mark Tuttle, and Lenore Zuck, for their useful comments and criticisms; Johan van Benthem, Brian Chellas, David Makinson, and Krister Segerberg for their help in tracking down the history of modal logic; and T. C. Chen and Brian Coan for pointing out the quotations at the beginning of Chapters 2 and 3, respectively. Finally, the second and third authors would like to thank the students of CS 356 (at Stanford in the years 1984–1989, 1991–1992, and 1994), CS 2422S (at Toronto in 1990) and the course on Knowledge Theory (at the Weizmann Institute of Science in the years 1987–1995), who kept finding typographical errors and suggesting improvements to

the text (and wondering if the book would ever be completed), especially Gidi Avrahami, Ronen Brafman, Ed Brink, Alex Bronstein, Isis Caulder, Steve Cummings, John DiMarco, Kathleen Fisher, Steve Friedland, Tom Henzinger, David Karger, Steve Ketchpel, Orit Kislev, Christine Knight, Ronny Kohavi, Rick Kunin, Sherry Listgarten, Carlos Mendioroz, Andres Modet, Shahid Mujtaba, Gal Nachum, Leo Novik, Raymond Pang, Barney Pell, Sonne Preminger, Derek Proudian, Omer Reingold, Tselly Regev, Gil Roth, Steve Souder, Limor Tirosh-Pundak-Mintz, Maurits van der Veen, Orli Waarts, Scott Walker, and Liz Wolf.

# Preface to the Paperback Edition

Relatively few changes have been made for this edition of the book. For the most part, this involved correcting typos and minor errors and updating references. Perhaps the most significant change involved moving material from Chapter 7 on a notion called "nonexcluding contexts" back to Chapter 5, and reworking it. This material is now used in Chapter 6 to refine the analysis of the interaction between common knowledge and agreement protocols.

The effect of teaching a number of classes using the hardcover edition of the book can be seen in this edition. The second author would like to thank the students of CS 676 (at Cornell in the years 1996, 1998, and 2000) for their comments and suggestions, especially Wei Chen, Francis Chu, David Kempe, Yoram Minsky, Nat Miller, and Suman Ganguli. The third author would like to thank the students of the course "Knowledge and Games in Distributed Systems" (at the Technion EE dept. in the years 1998, 2000, and 2002) for their comments and suggestions, especially Tomer Koll, Liane Levin, and Alex Sprintson. We would also like to thank Jelle Gerbrandy for pointing a minor bug in Chapter 3, and Rohit Parikh for pointing out minor bugs in Chapters 1 and 2.

The second and third authors changed institutions between the hardcover and paperback editions. The fourth author moved shortly before the hardcover edition appeared. The second author is now at Cornell University, the third author is at the Technion, and the fourth author is at Rice University. We would like to thank these institutions for their support of the work on the paperback edition.