

CONTENTS

<i>Acknowledgments</i>	<i>xi</i>
Introduction	1
Chapter 1: What Is Artificial Intelligence and Why Does It Raise Philosophically Interesting Ethical Issues?	6
1: The Blake Lemoine Incident	6
2: What Is Artificial Intelligence? An Evolving Concept	9
2.1: Thinking Machines? Turing vs. Jefferson	10
2.2. The Invention of the Term “Artificial Intelligence” and the Question of Whether AI Technologies Are Intelligent	13
3: What Definition of AI Should We Accept? Exclusive vs. Inclusive Approaches	17
4: Narrow and Broad Conceptions of—and Negative and Positive Approaches to—the Ethics of Artificial Intelligence	19
5: Why Does the Notion of Artificial Intelligence Raise Philosophically Interesting Ethical Questions?	21
6: Outlook	23
Chapter 2: Control and Value Alignment Problems, AI Hype, and the Notion of Artificial General Intelligence	27
1: The 2022 Release of ChatGPT and the AI Hype and Worries that Followed	27
2: What Are the Control and Value Alignment Problems? And What Does “Artificial General Intelligence” Refer To?	31

3: Is There Really a Control Problem and, If So, (How) Could It Be Solved?	34
4: Human-Compatible AI and Superalignment	37
5: An Out-of-Control Hype that Distracts Our Attention Away from More Pressing Ethical Problems Associated with Contemporary AI?	42
6: Do We Have to Make a Choice Between <i>Either</i> Focusing on Risks Related to Future AGI <i>or</i> Focusing on Harms and Risks Related to Today's AI Technologies?	46
7: Revisiting the Question of Whether AI Technologies Possess Any Forms of Agency and Intelligence	49
Chapter 3: Responsibility and AI, Good and Bad Outcomes, and Asymmetries Between Praise and Blame	57
1: Lee Sedol's Loss Against AlphaGo and the First Time a Self-Driving Car Killed a Human Being	57
2: A General Introduction to the Idea of Moral Responsibility	61
3: Ten Important Points About Moral Responsibility and About Praise and Blame and the Differences Between Them	64
4: Artificial Intelligence and Responsibility Gaps	70
5: Filling Responsibility Gaps by Having People Volunteer to Take Responsibility?	74
6: Praiseworthiness and Blameworthiness for Good and Bad Outcomes Created by/with AI Technologies	77
7: Concluding Remarks	81
Chapter 4: Generative AI's Gappiness with Respect to Meaning, Authorship, and Credit and Blame	84
1: From Hollywood Screenwriters to Poisonous Mushrooms	84
2: Gappiness?	87
3: Reminder: Credit-Blame Asymmetries	88

4: Meaning, Understanding, and Worries About a Lack of Meaning and Understanding in Generative AI Technologies	90
5: The Chinese Room and the Octopus	94
6: The Death of the Author, the Return of the Author, and the Question of Whether Generative AI Technologies Can Be Authors	98
7: Authorship Gaps: Can Generative AI Be an Author or Co-Author?	102
8: Should Human Users Always Be Considered the Authors of the Outputs of Generative AI Technologies?	106
9: Do We Need New Norms, Concepts, and Theories to Deal with Generative AI's Gappiness?	107
Chapter 5: Artificial Intelligence and Human Enhancement: Do AI Technologies Make Us More or Less (Artificially) Intelligent?	112
1: Back to the Man Who Was Carrying Out the Winning Moves in the AlphaGo Games Against Lee Sedol	112
2: Human Intelligence and Cognitive Enhancement	115
3: Artificial Intelligence and Its Relation to Human Intelligence	118
4: Might AI Technologies Enhance Human Intelligence by Extending Our Minds?	120
5: Artificial Human Intelligence	125
6: Concluding Discussion	129
Chapter 6: Artificial Intelligence, Consciousness, and the Moral Status of Conscious and Non-Conscious AI Agents	138
1: People in the Tech World Who Claim that AI Technologies Are Conscious	138
2: The Basic Concepts in This Chapter	141
3: A Quick Way of Mapping Some of the Relevant Recent Literature About AI Consciousness (or Lack Thereof) and Moral Status (or Lack Thereof)	147

4: Are Any AI Agents Conscious and/or Sentient or Might They Soon Be?	149
5: Consciousness and Moral Status	153
6: AI Agents Without Consciousness/Sentience but with Moral Status?	156
7: Critical Discussion of the Views Presented in the Previous Section	159
Chapter 7: Personalized Imitation Games: Permissibility, Desirability, and the Ethics of Digital Duplicates and AI Copies of Real People	168
1: From AI-Generated Pop Stars to Digital Duplicates of Dead Relatives	168
2: What Is a Personalized Digital Duplicate/AI Imitation of a Particular Person?	171
3: Personalized Imitation Games and AI Replacement Worries	174
4: The Permissibility of Personalized Digital Duplicates	178
5: The Desirability of Personalized Digital Duplicates	184
6: Conclusions of this Chapter and Some Further Discussion of Digital Duplicates	189
Chapter 8: Moral Principles for Human-Human Interaction, Human-AI Interaction, and AI-AI Interaction	195
1: Should AI Technologies Follow the Same Moral Principles as (Particular) Human Beings?	195
2: Should AI Technologies Follow Moral Principles in the First Place?	198
3: Moral Agents and Moral Patients, and Three Kinds of Interaction: Human-Human, Human-AI, and AI-AI Interaction	202
4: Traditional Ethical Principles and Principles for the Ethics of AI	208
5: The Double Standards Argument	211
6: Possible Justifications for Both Double and Triple Standards	214
7: Concluding Remarks and Transition to the Next Chapter	218

Chapter 9: The Future of (Meaningful) Work, Meaning in Life, and Human Flourishing in the Age of AI	223
1: Benefits for All of Humanity?	223
2: Different Kinds of Improvements to Human Lives	225
3: Can AI Make Human Lives Morally Better?	228
4: Five Quick Arguments for Why AI Can Have an Impact on Meaning in Life	231
5: The Achievement Gap and the Future of (Meaningful) Work	235
6: The Prospects for Meaning and Human Flourishing Within Futuristic AI Utopias	240
7: The AI and Meaning Sweet Spot	246
Chapter 10: AI Risks and Opportunities, Sustainability, and the Ethics of Gambling with the Future of Humanity	252
1: The AI Wager	252
2: The Ethics of Risk Imposition	255
3: AI and Different Kinds of Sustainability	260
4: The Drunk Driver Analogy and the “Who Are the Moral Agents and Who Are the Moral Patients?” Problem	264
5: Utopian Optimism vs. Dystopian Pessimism, and Responsible and Reckless Types of AI Optimism	268
6: A Young, Developing, and Quickly Growing Field	270
Concluding Summary: Retracing Our Steps	275
<i>Bibliography</i>	280
<i>Index</i>	294

INTRODUCTION

What is artificial intelligence (AI), and why does it raise philosophically interesting ethical questions? Why is it important to control AI, and what human values should AI technologies be aligned with? Who is responsible—who can be praised or blamed—when AI technologies produce good or bad outcomes? When people use generative AI to produce texts, images, music, or other outputs, who should be considered the author(s) of those outputs? If AI technologies continue to become increasingly “smarter,” is there a risk that human beings will become “dumber” if we begin to rely very heavily on artificial intelligence in all areas of life? What if it becomes possible to create AI with some form of consciousness? What will the future of (meaningful) work look like in a world in which we hand over more and more tasks to AI technologies? Does having more AI in all areas of life necessarily mean a better world for all of us?

These kinds of questions are discussed within philosophical debates about the ethics of AI. This book is about such questions. It explores different ways we can reflect on—and begin to answer—these and related questions from a philosophical point of view. The book is not directly concerned with concrete policy making, nor with legal or economic questions about our present and future life with AI. Nor is it concerned with sociological or psychological questions. It is about distinctively philosophical questions about the ethics of AI.

What does this mean? What makes something a distinctively philosophical question about a topic (such as the ethics of AI)?

Philosophers sometimes shy away from saying what exactly philosophy is. It can be many different things. Yet, there are some things that many philosophical discussions have in common. Here are some examples.

Philosophical discussions are often about basic concepts and basic ideas we have about a subject matter. A typical philosophical question about AI, for instance, could regard what exactly we should understand by the term

“intelligence” and whether it is possible to create technologies with something that we could call intelligence.

Philosophical questions are also often about basic principles or fundamental norms related to different parts of life. Which basic ethical principles, for instance, should we follow in life in general? And which ethical principles should govern our creation, deployment, and use of AI technologies?

Philosophical questions are sometimes also about apparent conflicts among different beliefs and convictions that many people have. For example, many people believe that we have a free will that is not determined from the outside. At the same time, they might also believe that everything that happens in the universe is determined by natural laws. Is this a clash of ideas or a contradiction, or could both things be true at the same time?

Philosophical inquiry is also sometimes about different ways in which things could be or not be in hypothetical scenarios (in different “possible worlds”). For instance, is it conceivable that there is a possible world where human beings have brains and nervous systems just like ours, but where these alternative human beings are complete “zombies” in the sense that they lack any subjective consciousness? Would we have different ethical duties toward such philosophical zombies than those that we have toward normal human beings? Or, to give an example related to AI: Could AI technologies that function in the ways that current AI technologies do become conscious, and if so, what would that mean for how we should interact with them?

A final example of a classic philosophical topic is what it is to live a good and meaningful human life. In the context of AI, we can ask whether AI technologies might create new ways in which human beings could live good or meaningful lives, or whether more and more AI in all parts of life might, in certain ways, make it harder for us to live good and meaningful lives.

This book relates these types of philosophical questions to the ethics of AI. It encourages you (the reader) to make up your own mind about these topics. And it will try to provoke you to reflect deeply on—and perhaps even worry about—these topics.

The book is not a neutral introduction to philosophical aspects of the ethics of AI. It is an opinionated introduction. I have chosen topics I find particularly interesting. And my discussion is partly driven by my own

philosophical convictions in what follows. The aim of making this an opinionated introduction is not to convince all readers that they should necessarily share all my views. Instead, the idea is that a book that offers arguments and critically assesses others' ideas is more interesting and more engaging than a book that tries to remain completely neutral.

The reader is encouraged to question and disagree with the ideas and arguments discussed in these chapters. That is what philosophy is all about. It is about mutually respectful and careful critical reflection—which sometimes involves deep disagreements—about the most basic questions of human life.

* * *

The book will consider three different kinds of questions:

1. Ethical questions related to actual AI technologies that already exist in society
2. Ethical questions raised by how people perceive AI technologies, or the assumptions people make about AI, whether or not these are completely accurate or well-founded
3. Ethical questions raised by AI technologies that do not yet exist, but that may come into existence in the near or distant future

The first two types of questions will be discussed at greater length than the third type. But the third type of question will be discussed as well. Such questions must be discussed with care. Otherwise, philosophical discussions about them may veer off too far in the direction of science fiction. At the same time, we should not forget that good science fiction can be very philosophical. Reflection on what is, at present, best considered as science fiction can be interesting from a philosophical and ethical point of view. Still, most of the book will be about real people, real events, actually existing technologies, and AI of the sorts we already have in society. But we will not shy away from occasionally reflecting on philosophically interesting questions about possible future forms of AI.

There are some topics often discussed within the ethics of AI—such as explainability, biases in data sets and outputs from AI technologies, or the environmental impact of AI technologies—that we will address to some

extent, but less than might be found in other books about the ethics of AI. There are so many different interesting and important topics we could discuss in a book like this. We cannot fit them all into one book, so I have had to make choices about what to focus on.

Speaking about making choices, in this context, we also have to make decisions about what types of AI technologies to discuss. When I first started researching the ethics of AI a decade ago, the hottest topic was the ethics of self-driving cars. Now, at the time of writing, the hottest topic within the ethics of AI is generative AI, and large language models in particular. A decade from now—or even sooner than that—the hottest topic might be something else. For instance, some people have recently suggested that the next big thing in the ethics of AI is the topic of “advanced AI assistants,” sometimes simply referred to as “agents” in the tech industry. At any rate, “agentic AI” is the latest hype in the tech world. So, how do we choose what AI technologies to discuss in a book like this?

One thing to notice is that a lot of philosophical questions about different AI technologies are actually similar in nature. In other words, the same questions tend to arise regarding different forms of AI within the ethics of AI at different points in time. For instance, will we lose control over these forms of AI? And who is responsible if something goes wrong? Those questions always come up. Then there are some questions specific to—or at least more clearly relevant in relation to—specific forms of AI. For instance, there are questions about language that are particularly relevant to current forms of large language model-based generative AI technologies. Such questions might be less relevant to, say, self-driving cars. This book features a mix of ethical questions that come up in relation to all forms of AI technologies, and ethical questions specifically related to particular types of AI.

From Chapter 1 onward, I will seldom mention myself. But as noted above, this is an opinionated introduction to the ethics of AI. So, what follows reflects the views and philosophical interests of the book’s author. I am passionate about and deeply fascinated by the topics I cover here. My hope is that readers will be inspired to think more about these and related topics, and perhaps themselves go on to write about them.

This book is partly meant to be used in introductory courses about the ethics of AI at universities. But it is also intended to be of interest to fellow

philosophy researchers who work on the ethics of AI, or anybody interested in the topic. Accordingly, I am assuming that you, the reader, are interested in artificial intelligence and philosophy, including the branch of philosophy that deals with ethics. So, if you are ready for some ethics of AI from a philosophical perspective, let's dive right in!