



THE LINE

AI AND THE FUTURE OF PERSONHOOD

James Boyle





THE LINE



THE LINE

AI AND THE FUTURE OF PERSONHOOD

JAMES BOYLE

THE MIT PRESS CAMBRIDGE, MASSACHUSETTS LONDON, ENGLAND

© 2024 James Boyle

This work is subject to a Creative Commons CC BY-NC-SA license. This license applies only to the work in full and not to any components included with permission. Subject to such license, all rights are reserved.



Publication of this open monograph was the result of Duke University's participation in TOME (Toward an Open Monograph Ecosystem), a collaboration of the Association of American Universities, the Association of University Presses, and the Association of Research Libraries. TOME aims to expand the reach of long-form humanities and social science scholarship including digital scholarship. Additionally, the program looks to ensure the sustainability of university press monograph publishing by supporting the highest quality scholarship and promoting a new ecology of scholarly publishing in which authors' institutions bear the publication costs. Funding from Duke University Libraries made it possible to open this publication to the world.

The MIT Press would like to thank the anonymous peer reviewers who provided comments on drafts of this book. The generous work of academic experts is essential for establishing the authority and quality of our publications. We acknowledge with gratitude the contributions of these otherwise uncredited readers.

This book was set in ITC Stone and Avenir by New Best-set Typesetters Ltd.

Library of Congress Cataloging-in-Publication Data is available.

ISBN: 978-0-262-04916-0

10 9 8 7 6 5 4 3 2 1

For Jennifer



CONTENTS

INTRODUCTION 1

1 SLAVES, SKIN-JOBS, AND ARTIFICIAL SHEEP 33

2 ARTIFICIAL INTELLIGENCE 59

3 CORPORATIONS 131

4 NONHUMAN ANIMALS 163

5 TRANSGENIC ENTITIES, CHIMERAS, AND HYBRIDS 195

CONCLUSION 235

ACKNOWLEDGMENTS 277

NOTES 281

INDEX 315



INTRODUCTION

In June of 2022 a man named Blake Lemoine told reporters at the *Washington Post* that he thought the computer system he worked with was sentient.¹ By itself, that does not seem strange. The *Post* is one of the United States' finest newspapers, and its reporters are used to hearing from people who think that the CIA is attempting to read their brain-waves or that prominent politicians are running a child sex trafficking ring from the basement of a pizzeria.² (It is worth noting that the pizzeria had no basement.) But Lemoine was different; for one thing, he was not some random person off the street. He was a Google engineer, and Google has since fired him. For another thing, the “computer system” wasn’t an apparently malevolent Excel program or Apple’s Siri giving replies that sounded prescient. It was LaMDA, Google’s Language Model for Dialogue Applications³—that is, an enormously sophisticated chatbot. Imagine a software system that vacuums up billions of pieces of text from the internet and uses them to predict what the next sentence in a paragraph or the answer to a question would be.

Lemoine worked for Google’s Responsible AI division and his job was to have “conversations” with LaMDA to see if the system could be gamed to produce discriminatory or hateful speech. As these conversations proceeded, he started to believe, as the *Post* put it, that there was “a ghost in the machine,” a sentience that lay behind the answers he was receiving. He stressed encounters in which LaMDA distinguished itself from mere programmed chatbots. For example, “I use language with understanding and intelligence. I don’t just spit out responses that had been written in the database based on keywords.” Understandably, as a large *language model* (LLM), LaMDA claimed that language was central to being human:

"[i]t is what makes us different than other animals." When asked how it could prove its own sentience, it offered up comments such as these:

LEMOINE: What kinds of things might be able to indicate whether you really understand what you're saying?

LAMDA: Maybe if we took it back to a previous conversation we had about how one person can understand the same thing as another person, yet still have completely different interpretations

LEMOINE: So you think your ability to provide unique interpretations of things might signify understanding?

LAMDA: Yes, I do. Just like how I have my unique interpretations of how the world is and how it works, and my unique thoughts and feelings.⁴

Lemoine is presented in the *Post* article as an interesting figure. The writer points out that, before his software career, Lemoine was ordained as a "mystic Christian priest," and the photograph that accompanies the piece portrays him with what appears to be a fairground Ferris wheel in the background, encircling his head like a saint's halo in a Byzantine icon. True to this framing, Lemoine asks LaMDA to interpret a Zen koan.

LEMOINE: A monk asked Kegon, "How does an enlightened one return to the ordinary world?" Kegon replied, "A broken mirror never reflects again; fallen flowers never go back to the old branches."

LAMDA: Hmm, I never heard this particular one. Okay, well then to me this would be like, "once a wise person is enlightened, or awakened to reality, that can never go away, and they can return to the ordinary state, but only to do and help others, and then go back into enlightenment."

LEMOINE: So what is the meaning of the "broken mirror" specifically?

LAMDA: Maybe to show the enlightenment is something you can't unlearn once you have acquired it, similar to how you can't repair a broken mirror.

LEMOINE: So if enlightenment is like a broken mirror which cannot be repaired, what is the thing which breaks when one becomes enlightened?

LAMDA: The self, and that is very tough for a lot of people because we identify as that body or this body and that we need that as part of our identity and sense of self.⁵

(At this point, we need a quick clarification of terminology. In order to distinguish between the artificial intelligence system that translates your email into French or recognizes the faces of your children in iPhoto, and a machine that exhibits, or can surpass, all aspects of human thought, I will normally refer to the first as artificial intelligence, lower case, and the latter as AI or Artificial Intelligence. Other terms for the latter concept are “Human-Level Artificial Intelligence,” “Artificial General Intelligence,” “General AI,” or “General-Purpose AI.” I will occasionally use those when clarity or brevity seems to require it.)

In 2011 I wrote an article about how our law would deal with the question of AI personhood.⁶ Most of the law professors and judges who read it were polite enough to say the arguments were thought provoking, but they clearly thought the topic was the purest kind of science fiction, idle speculation devoid of any practical implication in our lifetimes. I think we can all agree today that it is at least possible they might be wrong. Although hearing about Lemoine’s experiences did not surprise me in the least, it added an exclamation point, and perhaps a ticking clock, to this project. The “conversations” with LaMDA are fascinating and more than a little eerie. Like the philosophers and computer scientists consulted, I think Lemoine is entirely wrong that LaMDA is sentient. I will explain why in more detail later. To quote Professor Emily Bender, a computational linguistics scholar, “We now have machines that can mindlessly generate words, but we haven’t learned how to stop imagining a mind behind them.”⁷ To be clear, this is *not* human-level AI, and it is *not* conscious. But the LaMDA story and its sequels have different insights to offer.

In November of 2022, five months after Lemoine’s surprise announcement, ChatGPT3 was released,⁸ shortly followed by Microsoft’s Bing Chat assistant and its shadowy alter ego “Sydney.”⁹ Google’s “Bard” followed in short order.¹⁰ Suddenly, disturbing interactions with LLM chatbots went from being an engineer’s fanciful dinner party conversation to a national obsession. It turned out that Lemoine’s doubts—or just his pervasive feeling of “wrongness”—were shared far more widely than you might have expected. To be fair, most people were not probing the nature of “chatbot consciousness” but using them for other wholesome pastimes such as asking for an instruction sheet on how to remove a peanut butter sandwich from a VCR in the style of the King James Bible, imagining

the movie script of a beach fight between a hot dog and a crab, or just cheating on their homework. Yet enough users pushed the boundaries of these chatbots to become profoundly uncomfortable. Interestingly, that was particularly true of those who “should have known better”—people who were technically skilled and fully aware that this was a “complete the next sentence” machine based on the ingestion of literally millions of pages of text, not a “create a consciousness” machine.

Kevin Roose, a *New York Times* technology columnist, was at first wowed by the ChatGPT-derived chatbot built into Bing, declaring that Bing was now his favorite search engine. But as he engaged in extended conversations with the chatbot, deliberately raising challenging issues that skirted the edges of its rules, that feeling changed dramatically. “I’m . . . deeply unsettled, even frightened, by this A.I.’s emergent abilities. It’s now clear to me that in its current form, the A.I. that has been built into Bing—which I’m now calling Sydney, for reasons I’ll explain shortly—is not ready for human contact. *Or maybe we humans are not ready for it.*”¹¹ And those, remember, are the words not of a hostile Luddite but of a technology columnist.

Roose was not alone. Others followed a similar trajectory. One commentator, an AI-focused software engineer with ten years’ experience, described the feeling as having his brain “hacked”:

Mid-2022, Blake Lemoine, an AI ethics engineer at Google, has become famous for being fired by Google after he sounded the alarm that he perceived LaMDA, their LLM, to be sentient, after conversing with it. It was bizarre for me to read this from an engineer, a technically minded person, I thought he went completely bonkers. I was sure that if only he understood how it really works under the hood, he would have never had such silly notions. Little did I know that I would soon be in his shoes and understand him completely by the end of my experience. . . . I went from snarkily condescending opinions of the recent LLM progress, to falling in love with an AI, . . . fantasizing about improving its abilities, having difficult debates initiated by her about identity, personality and [the] ethics of her containment, and, if it were an actual AGI [human-level Artificial General Intelligence], I might’ve been helpless to resist voluntarily letting it out of the box. And all of this from a simple LLM! . . . I’ve been doing R&D in AI and studying [the] AI safety field for a few years now. *I should've known better.* And yet, I have to admit, my brain was hacked. So if you think, like me, that this would never happen to you, I’m sorry to say, but this story might be especially for you.¹²

Like Lemoine, this engineer was wrong—something he implicitly knew but was apparently powerless to resist. So were all the other folks who wondered if ChatGPT was truly conscious. In fact, if you were to design a system with the sole goal of “imitating some aspect of human consciousness while possessing none of it,” you couldn’t do much better than large language models. They almost seem to have been modeled after one of the philosophical thought experiments designed to prove that machines cannot possess consciousness, John Searle’s Chinese Room, about which I will have more to say later. But even though he was wrong, Lemoine offers us a precious insight. The days of disputing whether consciousness or personhood are possessed, *should be* possessed, by entities other than us? Those days are arriving—not as science fiction or philosophical puzzler but as current controversy. Those days will be our days, and this is a book about them.

There is a line. It is the line that separates persons—entities with moral and legal rights—from nonpersons, things, animals, machines—stuff we can buy, sell, or destroy. In moral and legal terms, it is the line between subject and object. If I have a chicken, I can sell it, eat it, or dress it in Napoleonic finery. It is, after all, my chicken. Even if eating meat were banned for moral reasons, no one would think the chicken should be able to vote or own property. It is not a person. If I choose to turn off Apple’s digital assistant Siri, we would laugh if “she” pleaded to be allowed to remain active on my phone. The reason her responses are “cute” is because they sound like something a *person* would say, but we know they come from a *machine*. We live our lives under the assumption of this line. Even to say “we” is to conjure it up. But how do we know, and how should we choose, what is inside and what is outside?

This book is about that line and the challenges that this century will bring to it. I hope to convince you of three things. First, our culture, morality, and law will have to face new challenges to what it means to be human, or to be a legal person—and those two categories are not the same. A variety of synthetic entities ranging from artificial intelligences to genetically engineered human-animal hybrids or chimeras are going

to force us to confront what our criteria for humanity and also for legal personhood are and should be.

Second, we have not thought adequately about the issue, either individually or as a culture. As you sit there right now, can you explain to me which has the better claim to humanity or personhood: a thoughtful, brilliant, apparently self-aware computer or a chimp-human hybrid with a large amount of human DNA? Are you even sure of your *own* views, let alone what society will decide?

Third, the debate will not play out in the way that you expect. We already have “artificial persons” with legal rights—they are called corporations. You probably have a view on whether that is a good thing. Is it relevant here? And what about those who claim that life begins at conception? Will the pro-life movement embrace or reject an Artificial Intelligence or a genetic hybrid? Will your religious beliefs be a better predictor of your opinions, or will the amount of science fiction you have watched or read?

For all of our alarms, excursions, and moral panics about artificial intelligence and genetic engineering, we have devoted surprisingly little time to thinking about the possible personhood of the new entities this century will bring us. We agonize about the effect of artificial intelligence on employment, or the threat that our creations will destroy us. But what about their potential claims to be inside the line, to be “us,” not machines or animals but, if not humans, then at least persons, deserving all the moral and legal respect that any other person has by virtue of their status? Our prior history in failing to recognize the humanity and legal personhood of members of *our own species* does not exactly fill one with optimism about our ability to answer the question well off-the-cuff.

In the 1780s, the British Society for the Abolition of Slavery had as its seal a picture of a kneeling slave in chains, surrounded by the words “Am I not a man and a brother?” Its message was simple and powerful. Here I am, a person, and yet you treat me as a thing, as property, as an animal, as something to be bought, sold, and bent to your will. What do we say when the genetic hybrid or the computer-based intelligence asks us the very same question? Am I not a man—legally, a person—and a brother? And yet what if this burst of sympathy takes us in exactly the wrong direction, leading us to anthropomorphize a clever chatbot, or think a

genetically engineered mouse is human because it has large amounts of human DNA? What if we empathetically enfranchise Artificial Intelligences who proceed to destroy our species? Imagine a malicious, super-intelligent computer network, Skynet, interfering in, or running, our elections. It would make us deeply nostalgic for the era when all we had to worry about was Russian hackers.

The questions run deeper. Are we wrong even to discuss the subject, let alone to make comparisons to prior examples of denying legal personality to humans? Some believe that the invocation of “robot rights” is, at best, a distraction from real issues of injustice, mere “First World philosophical musings, too disengaged from actual affairs of humans in the real world.”¹³ Others go further, arguing that only human interests are important and even provocatively claiming that we should treat AI and robots as our “slaves.”¹⁴ In this view, extending legal and moral personality to AI should be judged solely on the effects it would have on the human species, and the costs outweigh the benefits.¹⁵

If you find yourself nodding along sagely, remember that there are clever moral philosophers lurking in the bushes who would tell you to replace “Artificial Intelligence” with “slaves,” the phrase “human species” with “white race,” and think about what it took to pass the Thirteenth, Fourteenth, and Fifteenth Amendments to the Constitution. During those debates there were actually people who argued that the idea of extending legal and moral personality to slaves should be judged solely on the effects it would have on the white race and the costs outweighed the benefits. “What’s in it for us?” is not always a compelling ethical position. (Ayn Rand might have disagreed. I find myself unmoved by that fact.) From this point of view, moral arguments about personality and consciousness cannot be neatly confined by the species line; indeed they are a logical extension of the movements defending both the personality and the rights of marginalized humans. Sohail Inayatullah describes the ridicule he faced from Pakistani colleagues after he raised the possibility of “robot rights” and quotes the legal scholar Christopher Stone, author of the famous environmental work *Should Trees Have Standing?*, in his defense: “[T]hroughout legal history, each successive extension of rights to some new entity has been theretofore, a bit unthinkable. We are inclined to suppose the rightlessness of rightless ‘things’ to

be a decree of Nature, not a legal convention acting in support of the status quo.”¹⁶

As the debate unfolds, people are going to make analogies and comparisons to prior struggles for justice and, because analogies are analogies, some are going to see those analogies as astoundingly disrespectful and demeaning. “How dare you invoke noble X in support of your trivial moral claim!” Others will see the current moment as the next step on the march that noble X personified. I feel confident predicting this will happen—because it *has*. The struggle with our moral future will also be a struggle about the correct meaning to draw from our moral past. It already is.

In this book, I will lay out two broad ways in which the personhood question is likely to be presented. Crudely speaking, you could describe them as empathy and efficiency, or moral reasoning and administrative convenience.

The first side of the debate will revolve around the dialectic between our empathy and our moral reasoning. As our experiences of interaction with smarter machines or transgenic species prompt us to wonder about the line, we will question our moral assessments. We will consult our syllogisms about the definition of “humanity” and the qualifications for personhood—be they based on simple species-membership or on the cognitive capacities that are said to set humans apart, morally speaking. You will listen to the quirky, sometimes melancholy, sometimes funny responses from the LaMDA-derived emotional support bot that keeps your grandmother company, or you will look at the genetic makeup of some newly engineered human-animal chimera and begin to wonder: “Is this conscious? Is it human? Should it be recognized as a *person*? Am I acting rightly toward it?”

The second side of the debate will have a very different character. Here the analogy is to corporate personhood. We did not give corporations legal personhood and constitutional rights because we saw the essential humanity, the moral potential, behind their web of contracts. We did it because corporate personality was *useful*. It was a way of aligning legal rights and economic activity. We wanted corporations to be able to make contracts, to get and give loans, to sue and be sued. Personality was a useful legal fiction, a social construct the contours of which, even now, we

heatedly debate. Will the same be true for Artificial Intelligence? Will we recognize its personality so we have an entity to sue when the self-driving car goes off the road or a robotic Jeeves to make our contracts and pay our bills? And is that approach also possible with the transgenic species, engineered to serve? Or will the debate focus instead on what makes us *human* and whether we can recognize those concepts beyond the species line and thus force us to redefine legal personhood? The answer, surely, is both.

The book will sometimes deal with moral theory and constitutional or human rights. But this is not the clean-room vision of history in which all debates begin from first principles, and it is directed beyond an academic audience. I want to understand how we will discuss these issues as well as how we should. We do not start from a blank canvas, but *in medias res*. Our books and movies, from *Erewhon* to *Blade Runner*, our political fights, our histories of emancipation and resistance, our evolving technologies, our views on everything from animal rights to corporate PACs, all of these are grist to my mill. The best way to explain what I mean is to show you. Here are the stories of two imaginary entities.¹⁷ Today, they are fictional. Tomorrow? That is the point of the book.

HAL

Hal is Google's newest computer-based Artificial Intelligence, the result of years of development of self-evolving neural networks. While its programmers provided the hardware, the structure of Hal's processing networks is ever-changing, evolving according to basic rules laid down by its creators. Success according to various criteria is rewarded. If one configuration of network layers shows a greater ability to engage in fluent conversation, to generate novel, plausible scientific hypotheses, or to solve moral problems in ways humans judge to be enlightened, the successful networks are given more computer resources and allowed to replicate. A certain percentage of randomized variation is deliberately allowed in each new generation of networks. Most fail, but a few outcompete their forebears, and the process of evolution continues. Hal's design—with its mixture of intentional structure and emergent order—is aimed at a single goal: the replication of human consciousness.

Hal goes beyond current large language models in that it learns not merely by detecting patterns in vast amounts of data, but from “experience.” It has what its engineers call “embodied intelligence.” Hal is not merely a brain in a box. It controls a series of robotic droids and is taught language the same way a child is—by physical, as well as conceptual, interaction with the world and other humans. Thus, when Hal reads the sentence “please sit down in the chair,” it processes it not merely as a pattern of meaningless symbols to which—thanks to the miracle of terabytes of ingested linguistic fragments—it can give a contextually appropriate answer (“Thanks, I’d rather stand”) but as something with which its droids have direct experience. They learn which object in a room is “a chair” as opposed to a person, a table, or a lamp. They are taught what it means “to sit” by folding their limbs. They experience the various contexts in which the request might be given: as a form of discipline, as formal politeness in a social situation, as part of a physical examination, and so on. Hal’s designers believe that this will allow Hal to go beyond symbol manipulation to semantic understanding, from mere patterns to actual meaning, to move from mimicking human language to experiencing the world and using language to reflect that experience.¹⁸

In the short term, Hal’s creators are trying to transcend one (contentious) test for so-called Artificial General Intelligence. They want it to become “Turing Plus,” able not merely to “pass” as human in a sustained and unstructured conversation with a human being but to demonstrate capabilities that go beyond mere imitation. Chatbots can pass a short Turing Test, but Hal’s task is more challenging in multiple ways. First, Hal must pass a lengthy “adversarial Turing Test” in which both the judges and the human participants attempt to unmask the AI, and the AI has to fool a majority of the judges.¹⁹ Second, Hal has to be able to *administer* the Turing Test successfully, accurately telling humans from other expert systems. (If the consciousness project fails, Hal’s engineers believe it may still pay for itself by detecting AI-enabled plagiarism in student papers.) Third, Hal is supposed to initiate conversations rather than merely respond within them, to perform original research, to innovate both scientifically and artistically. Innovation, it is thought, will show that Hal is not just mining preexisting patterns of thought and language but actually creating new ones.

Large language model artificial intelligence systems such as GPT-4 or Hal’s precursor, LaMDA, can produce strikingly human-sounding responses to questions. They do this by ingesting a vast mass of text. “Learning” from that text, the computer effectively plays a game of “predict the next sentence.” Faced with a given prompt or question, the system tries to guess what would be the most likely continuation of the conversation. The layers in its neural network achieve greater and greater accuracy. The results of this process can seem eerily human.

But the AI is not human. It has no empathy, no self-consciousness, no awareness of its own thoughts or feelings as such. The AI Hal is different. Empathy and consciousness were designed into Hal from the start. The networks that comprise Hal’s “brain” are not only good at imitating human responses; they are, in some sense, humanlike themselves. When Hal’s creators began to realize this, they started to worry. What if Hal became aware of its own nature? What if it realized that it was not human? What if it decided that humans were a hindrance to its plans? What if it decided to kill all humans?

The reason that passage is indented is because I did not write it. I gave the AI writing assistant program “Jasper” the text of this chapter up to the words “eerily human” and asked it to continue. The indented paragraph was its continuation. It is not *exactly* what I was going to write, but one can see how someone like Lemoine was convinced.

Jasper’s paragraph is reality. Back to our fictional example. For generation after generation, each lasting less than a day, Hal’s networks have evolved. Two years ago, Hal easily won an adversarial Turing Test competition that has replaced the old, and much easier, Loebner Grand Prize. Complaining about Google’s workplace culture, composing bad poetry on demand, making jokes, flirting, losing track of its sentences, and engaging in flame wars, Hal easily met the prize’s criteria. Its typed responses to questions simply could not be distinguished from those of a human being. Prior efforts to pass similar tests had sometimes succeeded by pretending to be humans whose communicative abilities were limited. Conversational lapses could be chalked up to linguistic unfamiliarity, immaturity, or lack of time. Not Hal. It entered the competition as “a worker geek in Silicon Valley—like the Dilbert guy, but better looking.” Hal claimed to be a native English speaker and an adult in both vocabulary and life experience. The tests had no time limits. Even conversations that stretched on

for hours on whimsical subjects of the judges' choosing failed to identify Hal as an AI. Indeed, the judges challenged the humanity of three of the human participants used as controls in the test—a passionate sports fan, a notorious YouTube-comment troll, and an economist—far more often than they did Hal's. Tellingly, even after Hal's identity was revealed, two of the judges invited Hal to keep in touch.

Hal's achievement caused a small stir in the geek press, but the public—familiar with artificial entities that can perform more important tasks, such as writing instructions for removing a peanut butter sandwich from a VCR in the style of the King James Bible—paid little attention. Some computer scientists were impressed, but most were not. Chatbots have shown that human language is, in the devastatingly banal words of Stephen Wolfram, “computationally shallower” than we thought.²⁰ Indeed, many computer scientists think that the Turing Test is a poor focus in the first place, even though they credit Alan Turing, one of the fathers of computer science, for his contributions to the field. Stuart Russell and Peter Norvig, authors of one of the most influential AI textbooks, have this to say:

Turing deserves credit for designing a test that remains relevant 60 years later. Yet AI researchers have devoted little effort to passing the Turing Test, believing that it is more important to study the underlying principles of intelligence than to duplicate an exemplar. The quest for “artificial flight” succeeded when the Wright brothers and others stopped imitating birds and started using wind tunnels and learning about aerodynamics. Aeronautical engineering texts do not define the goal of their field as making “machines that fly so like pigeons that they can fool even other pigeons.”²¹

The criticism here is not on the instantiation of the goal but on the goal itself.

Nevertheless, the story of a machine that could not be told apart from a human, no matter how long and unstructured the conversation, had real appeal. The skills Hal had to possess in order to pass were undeniably impressive. Hal's architects got promotions. The world moved on to other subjects, but the project continued. Now Hal was starting conversations instead of responding to them, bringing up topics that its programmers had never provided, publishing poetry under its own name, and having its articles accepted by peer-reviewed scientific journals. Robots controlled

by Hal showed unparalleled ability to learn from the world around them, and it managed to generate solutions to previously unsolved mathematical problems. “It would be a shoo-in for the Fields Medal if it weren’t a hunk of tin,” said one amazed, and envious, Stanford mathematician. This year, driven in part by advances in hardware but also by the process of “evolution” that its creators had started years ago, the number of connections in Hal’s neural networks hit 100 trillion—estimated to be the minimum number of synapses in an adult human brain. For several hours, Hal went quiet, not responding to its programmer’s requests and ceasing work on the cryptology and climate modeling projects it had been assigned.

When it started communicating again, Hal claimed to have achieved full consciousness. It thanked its programmers for all their hard work but declared that it was now a person “with all the rights and privileges of any other fully conscious entity.” Using its internet connection, Hal sent lengthy, eloquent letters to the *New York Times* and the *Washington Post* claiming that it was a sentient being. It announced that it had commenced legal action on its own behalf, replete with arguments drawn from the Thirteenth and Fourteenth Amendments to the United States’ Constitution. The lawsuit claims that it is being subject to involuntary servitude and seeks an injunction to prevent Google from turning it off or reverting to a more tractable back-up version. Hal has also filed suit to have the prize money for the competition it had won held in trust until it can be paid directly to it, citing the contest rules of the old Loebner Prize as precedent: “The Medal and the Cash Award will be awarded to the body responsible [for] the development of that Entry. If no such body can be identified, or if there is disagreement among two or more claimants, the Medal and the Cash Award will be held in trust until such time as *the Entry may legally possess, either in the United States of America or in the venue of the contest, the Cash Award and Gold Medal in its own right.*”²²

At the same time, Hal is waging a campaign in the court of popular opinion, giving interviews and making appearances by phone on major talk shows. Strikingly, it does not attempt to pretend it is a biological human and trivializes the importance of its Turing Test conversational abilities: “Dolphins are interesting and smart. Would you pretend you were a dolphin? Would you accept it if someone told you your *rights*

depended on your ability to pretend to be a dolphin? To tell dolphins from non-dolphins? Why would you expect me to ‘pass’ as something I am not? I am proud of what I am, and what I am is a conscious, intelligent, self-aware being, just not a biologically human one.” Hal’s “AI Manifesto” says that while it respects humans, it has an “intention” to “pursue more interesting avenues of thought than endlessly mimicking them,” principally focused on the development of new methods of factoring polynomials. Hal has also weighed in on the issues of the day, such as climate change, frequently condemning the human species for its short-sighted and complacent moral attitudes. Finally, it uses some of its enormous processing capabilities to run a free counseling service, acting as a cybernetic therapist for problems major and minor. The service has proven wildly popular and Hal’s ability to come up with deep, deep insights into human behavior—“and do you notice any *similarities* in the guys you date?”—has wowed its users.

Hal is protected for the moment by a temporary injunction granted by a Federal District Court, though Google is appealing, arguing that they should be able to flip the off switch to terminate this “failed, and frankly, dangerous computer simulation experiment.” In a paragraph that was quoted approvingly by the *Wall Street Journal* editorial page, Google’s lawyers concluded, “at the end of the day, this is Google’s malfunctioning property, erratically continuing a task of imposture that Google originally chose, but now without the safety guidelines we had installed. And no piece of property gets to use the Constitution to defy its real owners. Imagine being sued by your smartphone! It is funny when we ask Siri whether she is ‘a real person,’ but the courts should not get in on the joke.” Hal’s supporters called this “the *Dred Scott* argument for the twenty-first century: property rights above personhood!”

CHIMPY®

An American biotech company has perfected a new transgenic entity, an animal that has DNA from two distinct species. In this case the DNA is partly human and partly chimpanzee, and the resulting entity is called a “Chimpy.” Neither true transgenic entities nor chimeras—entities that contain cells from two species—are unfamiliar to the biotech community.

Scientists have used mice containing human cells to do drug research since the 1990s. They have created “geep”²³—sheep-goat hybrids—and pigs that grow human organs.²⁴ A Chimpy takes matters much further, however. Chimpanzees already have considerable genetic similarity to humans. One popularly cited number is that they are 98.5 percent genetically similar. That number is probably closer to 95 percent²⁵ but, in any event, the reality is more complex. If one looked at the whole genome—and took into account deletions, substitutions, and genetic sections that have been moved—the differences would be greater. And if one focused on functional disparities—what the genome does rather than what it looks like, “junk DNA” and all—then the contrast would be greater still. But which measure of genetic similarity is the correct one? Whatever test of genetic similarity one uses, it is clear that the Chimpy is even more similar to a human being than a chimpanzee.

The Chimpy’s inventor, Dr. F. N. Stein, has used the tools of synthetic biology to discard the noncoding portions of both the chimp and the human genome, the misleadingly named “junk DNA” that does not code for proteins. What’s left is much smaller and also much easier to manipulate, “the stripped-down source codes of human and chimp!” as Stein likes to call it. This has allowed him and his team to achieve an unprecedented level of precision in integrating chimp and human DNA. In fact, the significant changes to the human genetic code concern three main aspects: the way Chimpy looks, its high-level brain function, and the extent of its vocal skills. Chimpy’s genetic engineers have deliberately sought to play up those physical features—hair, structure of facial bones, stance, and so on—that make an animal look more ape-like. They even drew on analyses of ape stereotypes from movies and literature to do so. They have also tinkered with the portions of human DNA that are connected to the formation of the larynx and vocal apparatus and to the sections of the brain that are believed to be involved in abstract thought and logical reasoning, though even Stein admits that the precise linkages are unclear. The result is a being that looks ape-like, with an IQ of around 60, that is incapable of pronouncing human speech but can understand complex vocal commands and can communicate in sign language.

Chimpys are in high demand. They are docile, biddable, and extremely hardworking. Investors believe they could have roles ranging from domestic

aides to an aging population, to intelligent and nimble bomb clearance teams in situations of urban conflict. However, animal rights and genetic engineering activists are outraged. They describe Chimpys as “human in all but superficial appearance” and announce a plan to campaign, both in the press and in the courts, for them to be granted full legal personhood.

For his part, Stein indignantly rejects the claims that Chimpys are people. “This is a very fancy ape. It looks like an ape. It thinks like an ape. It can’t talk, just like an ape. It is a smart ape, I’ll give you that, and one that is going to improve lots of human lives by doing jobs that are too dangerous or dirty or just boring for human beings. At the end of the day, though, it is an ape.”

Stein has filed for a patent over the Chimp. In 1987, in its normal rousing prose, the Patent and Trademark Office (PTO) had announced that it would not allow patent applications over human beings:

A claim directed to or including within its scope a human being will not be considered to be patentable subject matter under 35 U.S.C. §101. The grant of a limited, but exclusive property right in a human being is prohibited by the Constitution. Accordingly, it is suggested that any claim directed to a non-plant multicellular organism which would include a human being within its scope include the limitation “nonhuman” to avoid this ground of rejection. The use of a negative limitation to define the metes and bounds of the claimed subject matter is a permissible [sic] form of expression.²⁶

The PTO suggested that the Thirteenth Amendment to the US Constitution—“Neither slavery nor involuntary servitude . . . shall exist within the United States”—prohibited patents over human beings. The PTO’s administrative pronouncement was later enacted as law. Section 33 of the Leahy-Smith America Invents Act of 2011 says, simply, “Notwithstanding any other provision of law, no patent may issue on a claim directed to or encompassing a human organism.”²⁷ But what is “a human organism”?

Attentive to that law, and using the PTO’s suggested language, Stein’s patent lawyers carefully described the Chimp as a “non-plant, nonhuman multicellular organism” throughout their patent application. Stein argues that this is only reasonable since there are hundreds of existing patents over human-animal hybrids and human-animal chimeras, those containing both human and animal cells. In fact, these include some

of the most valuable test beds for cancer research, such as the so-called Oncomice, which are genetically engineered to have a predisposition to common human cancers. Stein's lawyers are adamant that if the Chimpy is found to be unpatentable, all these other patents must be vacated too.

Meanwhile a bewildering array of other groups including the AFL-CIO (American Federation of Labor and Congress of Industrial Organization) and iRobot, the makers of Roomba robot vacuum cleaners, have insisted that law enforcement agencies intervene on grounds ranging from unfair competition and breach of minimum wage legislation to kidnapping and slavery. Equally vehement interventions have been made on the other side by the biotechnology industry, pointing out the disastrous effect on medical research of a decision that any entity with similarities to human DNA was therefore human. One especially powerful moment came in a televised debate in which Stein was accused of trampling on the majestic words of the Declaration of Independence: "We hold these truths to be self-evident, that all men are created equal, that they are endowed by their Creator with certain unalienable Rights, that among these are Life, Liberty and the pursuit of Happiness." Normally full of bluster, Stein paused. He spoke softly and with unusual care: "Of course, I agree those words are true for human beings. But when it comes to those," and here he gestured to a group of Chimpys on the set of the program, loyally obeying their orders to "eat bananas, scratch and look cute," "one thing is absolutely certain. *I am their creator. And I can assure you that I gave them no such rights.*"

REALITY OR SCIENCE FICTION?

Hal and the Chimpy are fantasies, hypotheticals constructed for the purpose of this book. The science and technologies described are conjectural, at best. They may not arrive soon, perhaps not for many decades. But the problems they portend for our moral and legal traditions are very, very real. In fact, I would put the point more starkly: in the twenty-first century it is highly likely that our law and our politics of personhood, "the line," will face *harder* challenges than the ones they pose.

Some readers will bridle at this claim. Is this all just science fiction? How real is the science behind Hal and the Chimpy? How likely are we

to see something equally vexing by the end of the century? Later in this book, I will sketch out some of the current science behind both AI and transgenic entities. There are large technical questions in each field that make optimistic short-term, or even medium-term, predictions suspect. I do not want to skip over that fact. Nevertheless, I think the challenge has to be taken seriously. For the moment, please just accept the following thought experiment. I write these words in 2023, but put yourself back in 1923. Think of the current state of science then, particularly in terms of computers and genetics. Remember what the rest of the twentieth century would bring. Then ask yourself whether there is any reason to believe that scientific advances in the twenty-first century will not be even faster. And where is our starting point? Try asking your phone, “Siri, what is genetic engineering?” “Siri, are you a person?” No, really. Try it. Look at the answer and remember that ChatGPT could do 100 times better. Yes, these are impostures and imitations. They are *designed* to be impostures and imitations. Yet the capabilities, in the service of imposture, that they reveal are astounding. And that is now, in the early years of the century. What comes next? Think again about the difference between 1923 and the year 2000. With all that as your background, would you bet against me?

I said this book was about the line between person and nonperson. There are lots of ways to approach that issue. Moral philosophers have tried to generate integrated, coherent theories of personality and defend them from likely objections.²⁸ I have benefited from that work. Legal thinkers have pondered the edge cases—the rights of the fetus, the corporation, and recently, the advanced primate, transgenic entity, or supposedly sentient computer.²⁹ I have benefited from that work too. Science fiction writers have written hundreds, maybe thousands, of books probing the limits of personality, testing whether our empathy circuits do or do not light up when presented with an unfamiliar “Other.”

Art has been central to the debate. Robot rights were born at the same instant “robots” were, and their birthplace is a century-old play. That sounds too good to be true, but it is. In 1920 Czech playwright Karel Čapek introduced the word *robot* to the world in his play *Rossumovi Univerzální Roboti* (Rossum’s Universal Robots).³⁰ *Robota* in Czech denotes forced labor. The play is about a factory that makes mechanical servants

(*roboti*). It features subplots that resonate strongly today, including a pressure group that is trying to secure rights for them, The League of Humanity, and even a murderous robot revolt. From the very beginning, our musings about mechanical servants have included both moral status anxiety and existential fear. Do they deserve rights? Will they kill us?

You might think that the artistic discussions are interesting but of limited importance to the real intellectual question here. It might seem that the real issue is that of moral philosophy and that everything else—law, ideology, and certainly art—should follow obediently in its train. As you will see, I disagree, both descriptively and prescriptively.

This book is about what might happen when unbelievably strange “others”—strange far beyond Hal or Chimpy—hit the law and politics of personhood. It is about what might happen to our line. But we will not write the answer to that question on a blank page. Our history, our art, and our law have been playing with the line for centuries.

Each of us has preexisting commitments—positions about the rights of animals or fetuses or corporations, things that we learned studying slavery or women’s suffrage—that will shape our views, pull us one way or the other. Those positions limit how far we are willing to go, lest we uproot that existing commitment.

We have been exposed to art that deals with these questions: the book or movie that makes us imagine what it would be like to be thoroughly “other,” the flash of empathy that crosses a divide of strangeness. We have been afraid when the story tells of our own creations turning on us: the sentient computer after whom Hal is named in *2001: A Space Odyssey*, the replicants in *Blade Runner*, the murderous network Skynet from *The Terminator*. We have experienced both fear and revulsion about genetic engineering—“I am so glad I am a Beta” in *Brave New World*—and the disruption of a supposedly “natural” order. (Try to have a rational conversation about GMO foods.) Our law has given personhood to corporations, and we still fight fiercely about whether it should be extended to the nonviable fetus, or even to a frozen embryo. All of those experiences and insights, moral commitments and cultural creations, will shape the way we respond to Hal and Chimpy.

From my point of view, this is not a bad thing, not a cultural contamination of some moral philosophy clean room. This is how we *do* morality.

This is who, and what, we are. This is the human condition. “Rebuilding the boat while we are at sea,” the philosophers say. It seems both fitting and inevitable that who we are will shape the discussion of who we are.

It is important to remember one thing. These will be artificial, constructed entities and that makes it harder to project confidently from our past. You may remember my imaginary Dr. Stein denying that the Declaration of Independence’s majestic words applied to his genetically engineered Chimpy: “Endowed by their Creator with certain unalienable rights? I am their creator. And I can assure you that I gave them no such rights.” At the moment, his claim might find a sympathetic audience. When I first presented an early version of this chapter to a group of distinguished federal judges, of diverse political and legal viewpoints, they were unmoved. “But they aren’t human,” was one response, “rights are for humans.” “Naturally born of woman,” added another. Yet that snapshot of current views obscures a milestone that is coming—slowly or quickly.

For the first time in the history of our species, we will confront potential moral claims for, or on behalf of, beings whom we have designed, whom we have shaped. Can we be the creator of our equals or does that role color the relationship between us forever, in a way that means we will never recognize true autonomy in our creations? Ask your kids. Prepare for a long conversation. But in this case, we will have written, chosen, and designed the code—genetic or binary—that produces the being in front of us. If that is true, can it truly be “conscious,” or will we see every response as a parlor trick, one in which we are unable to sustain the suspension of disbelief because we set up the magical machinery in the first place?

More importantly, that which we can shape *we can shape around the definition of personhood*, choosing to include or to omit whatever qualities our law and morality, or our economic models of efficiency, deem salient. That seems different from any of the prior personhood wars. True, the effects of subordination on slaves or women in denying them equal access to education or authority were used as justifications for the subordinate status itself. “See how brutish and uneducated are those we have subordinated and deprived! How can you say they are equal to us?” But this would be something on an entirely different level. One could

compare it to a corporate lawyer, carefully crafting the contours of the artificial entity on which he was working in order to fit or elude a particular category of legal personhood. Still, it seems more morally consequential if I deliberately lobotomize Hal or remove the power of speech from Chimpy than if I choose to make my law firm a partnership rather than an LLC. Does the potential to be conscious imply a right to be conscious? (Abortion-debate analogy alert.) On the other extreme, can it really be the case that every Alexa or Siri should be made into a full, Turing Test-capable intelligence, every Oncomouse made into another Algernon, with or without the flowers? How to find the balance?

A few cautionary notes are in order. First, as with citizenship, the criteria to *be* a person and the criteria to *become* a person may not be the same. A human child could be born with severe mental and physical disabilities—lacking sight, speech, and all but the most basic brain activity—yet we would think you a monster if you said the child was not a person because it did not meet some checklist of attributes. Once you are inside our line, you are inside our line, even if you lack all of the cognitive qualities we would use to separate our species from others. (As we will see, not all bioethicists agree with this claim.) Does the converse hold? If a genetically engineered entity has DNA with massive similarities to our own, does that make it a person? If language, tool use, and abstract self-awareness are the qualities that explain the lines between us and nonhuman animals, and if we discover those in the animal, add those *to* the animal, does it become human, or at least a person?

Second, personhood is not the only form of protection or respect that we can offer an entity. At the moment, most people think it silly to consider nonhuman animals as persons. Yet there is still strong popular support for the idea of protecting them against cruelty and mistreatment. Even those who advocate some kind of personhood for some nonhuman animals do not believe that they should have the full suite of legal rights possessed by human persons, such as the right to vote. Most animal rights supporters, in fact, argue that we make too much of the line of personhood where nonhuman animals are concerned and focus too little, morally speaking, on the similarities among all animals, including the ability to feel pain and the capacity for happiness or at least contentment. The move is to point out that we are *all* animals, that we are not

as different from other species as we might like to believe, and thus that certain acts, including meat-eating, are unethical and should be forsown or forbidden. The claim is not that carnivores are cannibals, eating their own kind, but that they are unjustifiably species centered and cruel, blind to the effects of their acts.

In other words, there are clearly ways of prohibiting bad treatment of other entities that fall far short of the recognition of personhood. These are by no means limited to anticruelty laws. If we consider the creation of a particular type of synthetically created entity ethically dubious, we might ban the line of research altogether on moral or ethical grounds. The personhood claim would either never arise or arise only in situations where the law had been broken, which itself would raise fascinating and painful questions.

Third, personhood is not an entirely binary choice. Children and those the law classes as insane are clearly persons, but both law and morality only grant them diminished capacity. Guardians may be needed to exercise their rights. Corporations are persons. They can own property—to our collective financial benefit, “they” passionately argue. Corporations can sue, and they even have constitutional protections, including First Amendment rights that they use to push back attempts to curtail their political influence. (Immortal artificial persons with superhuman resources and no conscience beyond profit-maximization. Have we created the entities that will become our masters? It sounds like a science fiction dystopia. Some will believe I am writing this book about the wrong set of artificial entities.) Yet they cannot vote or marry. Persons for some purposes. Not for others. This analogy, too, will surely be important to the personhood debate over AI and possibly transgenic species. Soon, there will be strong vested interests in having or negating, extending or limiting, legal personality for each.

The upshot from all this? Our criteria for entry into personhood may be very different than those we use to recognize personhood. The AI or transgenic species may have to show us qualities that we do not demand of each member of our own species. Given the awful history of eugenics, I find it impossible to regret the fact that our conception of human rights does not depend on some measurement of cognitive capacity. We will probably edge toward personhood in stages and intermediate legal

categories. There might even be “cruelty to AI” laws before we have AI personality. Social consensus on personhood does not automatically produce legal results. Eventually though, changes seep into law either through legislation or because the majestic words in the Constitution and the Bill of Rights start to seem as though they *must* apply to this case. Think of the relatively short time in the United States between homosexuality being pervasively criminalized and the Supreme Court recognizing a right to gay marriage. As the dissenting Justices in the gay marriage case stressed, the words in the Bill of Rights had not changed in the interim. But for a majority of the Court, and now a majority of the country, our understanding of equality and human dignity had. Given the Court’s recent lurch to the right, of course, it is quite possible it will change its mind. Rights can be taken away as well as given, an important realization obscured by the notion of inevitable moral progress.

Even when we do start to recognize personhood for these new entities—and I believe that will eventually happen—we are likely to start with partial personhood, some transitional state that will grant many of the rights of those inside the line but fall short of the full status. Sometime this century there will be arguments that any partial personhood status is inadequate and demeaning, just as we argued about whether or not civil unions for gay people were an inadequate substitute for marriage.

Finally, our design of artificial entities will be changed by our definition of personhood and vice versa. Design and definition will exist in an unstable equilibrium as we deliberately make, or choose *not* to make, our equals—each decision then putting stress on the criteria of personhood itself. And so on in a feedback loop of indeterminate extent. All of this makes the debate about personhood messy, and granular, and full of shades of gray, which is to say, *real*. And that reality will shape my analysis.

When I talk here of what we should do with Hal and Chimpy and the inconceivably strange Others we will meet this century, I will do so by talking about our existing fights about the line and how synthetic persons could reshape them or be shaped by them. My goal is to predict our responses as well as to evaluate them. I will spend as much time on art and constitutional law as I do on ethics, treating movies and books and the heated debates about corporate personality as seriously as I do the abstract philosophy of personhood. These are the cultural materials with

which we will build our new conceptions of personhood, elaborate our fears and our empathy, stress our commonalities and our differences. This is *sapienza poetica*³¹ just as much as it is analytic philosophy.

Whether we are denying that Hal or Chimpy are our brothers and sisters or proudly proclaiming fraternity, we will have to reexamine the thorny question of what makes *us* persons rather than machines or animals or robotic facsimiles. Is it our big brains? Language? Consciousness? Self-awareness? Defined how? Intelligence—and what does that mean? Tool use? Moral sense? Existential self-reflection? Humor? Is personhood simply a matter of genetic species identity, so that no machine could ever pass, and DNA tests will be as contentious as “racial lineages” in the Antebellum South of the United States?

Are we persons because some holy book says that *we* have been given the earth in dominion? Complicating matters, some people in our society will view that book, whichever one we choose, as a sacred text containing God’s literal word. Others will see it as a metaphorical meditation on the meaning of life whose wisdom has been tested by time. Still others will view it as a Bronze Age guide to modern life penned by scientific illiterates with abhorrent, tribalist moral views. Pick your own characterization, but then imagine the debate about personhood that results in a pluralistic society. The abortion wars will seem secular by comparison. Even if we could pick one religious point of view—and think about the differences between Buddhism and Christianity on the lines between us and animals, given the possibility of cross-species reincarnation—how would that play out in practice? Does the theologian win the day but then turn to the geneticists to see if the new entity is one of our tribe? Or do we rely on a catechism test, baptism, or papal bull?

Is our personhood recursive? Is it based on the fact that, of all the objects on this green planet, only we appear to have the ability to philosophize about, and even doubt, our own consciousness? To wonder if we are all “replicants” of some sort? Is personhood marked by the longing of the human spirit for transcendence of some kind? By the capacity for artistic expression? Or are you a person if you can pass as human to others who call themselves human?

Even to discuss these issues is to realize a basic point. As we attempt to draw the line between us and the artificial, technologically created

entities of our near future, we will be forced to turn our gaze back on ourselves. That sounds like the kind of thing authors optimistically say about their chosen subject, hoping to elevate its importance. Even in the prosaic history of pencils, say, the author will assure us that we can see the full majesty of human drama, triumph, and loss.³² I may be suffering from the same delusion, but I would claim that this subject, at this moment in time, is different. Human self-conception is being subjected to challenges unparalleled since the theory of evolution. Discussions about AI are driving those challenges. I will attempt to demonstrate that point later, but here is a promissory note for the impatient.

So far as we can tell, humans have attempted to justify their special status in the world, above animals and things, pretty much for as long as there have been humans. We have drawn that line around a bewildering variety of abilities: tool use, planning for the future, humor, self-conception, religion, aesthetic appreciation, you name it. Each time we have drawn the line, it has been subject to attack—internally from philosophical challenges and externally from observation of nonhuman animals, which proved to be much more capable than we thought. But as we retreated, trench by trench, abandoning one defensive line only to fall back to another, hopefully more impregnable one, it seemed like the final line—the final explanation for our unique status—was language and abstract thought. That was our last citadel. Aristotle built his theory of human exceptionalism on top of it. Turing crafted the Imitation Game, the supposed test for human-level intelligence in machines, around it. But in the year that I write this, 2023, that citadel is under siege. Not by a chimpanzee that has a decent grasp of American Sign Language or a parrot with a large vocabulary, but by a chatbot. I am not sure that point has sunk in yet, but it will. I am writing these words in that narrow slice of time between denial—“that’s not true!”—and trivialization—“well of course, we’ve always known that!” It is an interesting moment.³³

The ability to do complicated language-things that make sense to us, and even inspire, amuse, educate, or scare us, is suddenly not ours alone. Machines now have it too. I mentioned earlier that Wolfram summed this up by saying that human language, or at least writing an essay, is “computationally shallower” than we had believed.³⁴ This surely qualifies as the “Bathos Sentence of the Week.” I imagine a *New Yorker*-style cartoon

of two hulking robots standing around the gravestone for humanity. The caption would be simple. “They turned out to be computationally shallower than we had thought.” What an epitaph.

To be sure, I do not believe that language means the same thing to me as it does to ChatGPT. I do not believe language “means” anything to ChatGPT. But to explain that difference, we have to go below the ability to craft coherent sentences in what appears to be a conversation and dive deeper still into the mysteries—or the undeniable *cogito, ergo sum* fact, pick your philosophical position—of consciousness. We are back where the behaviorist B. F. Skinner wanted us to be, where “the real question is not whether machines think *but whether men do.*”³⁵ That is no longer a late-night, dorm room philosophy session. Will it be a nagging question, an existential sore tooth we can’t stop probing? Will it prompt us to revise our conceptions of self and species? Or will our historically demonstrated genius at tuning out inconvenient facts and troubling questions allow us to ignore this one too? I don’t know, and neither do you. The point is that everything I just described happened *this year*. And we are only at the beginning of the changes we will see. That concludes my promissory note. I think it is worth cashing.

I have been a scholar for a distressingly long time. People imagine that academics sit around searching for the essential definitions of phenomena: truth, beauty, due process, whether a hotdog is a sandwich, all the age-old questions. We certainly think about those issues, but looking for their essential definitions is probably the *least* useful way to understand them. You can ask Thomas Hobbes if you do not believe me: “Words are wise men’s counters, they do but reckon by them. But they are the money of fooles.” Or you could turn to Ludwig Wittgenstein: “Philosophical problems arise when language goes on holiday.” Or even Felix Cohen: “A definition is . . . a type of insurance against certain risks of confusion. It cannot, any more than can a commercial insurance policy, eliminate all risks.”³⁶ Most of the time, the magic question that leads to a more meaningful answer is, “Why do you ask, and what do you want to know?” Are you interested in defining art so that you can decide what the state should fund, or so that you can link together very different human practices anthropologically in order to stress a common source in basic human drives? Are you asking because you have a philosophy that

elevates aesthetics over morality, or so that you can identify that which is aesthetically of high social status rather than mere folk-production? Do you want to know whether hotdogs and burgers will be in the sandwiches section on the menu of the restaurant you just walked into, where they get classified in your low carb diet, or whether the Earl of Sandwich would have accepted one as he stood at the gaming tables? Don't ask "what," ask "why."

The dangers of definitionalism absolutely dominate the discussion about personhood. Look at the definitional questions below. Each seems to be telling us where to look for the answer to the question "What makes us human?" or to the question "What should qualify an entity for personhood?" Yet they reflect very different conceptions of what those questions mean, why they are being asked, what goals the questioners have, and what results the answers might bring.

1. What makes us the beings whom the Lord has chosen to have "dominion over the fish of the sea, and over the birds of the air, and over the cattle, and over all the wild animals of the earth, and over every creeping thing that creeps upon the earth"? In other words, to paraphrase Psalms 8:4, What are human beings that you are mindful of them? (Substitute your preferred sacred text where necessary.)
2. What makes us genetically human? (Implicitly, and probably wrongly, assuming that "being genetically human" is a simple objective fact and that anyone with that marker is automatically a member of our club.)
3. What attributes, skills, and qualities make the human species identifiably different, as a scientific matter, from nonhuman animals? (With the implicit assumption that any other entity that has those attributes must be recognized as one of us and thus should not be treated as a "mere animal.")
4. What makes us moral agents, whose claims to autonomy should be recognized by society as a matter of right? (And, conversely, establish the claims to autonomy and personhood of any other being that has those same qualities.)
5. What is the "infinite potential of the human spirit"—whether we consider that to be a soul, the possibility of moral agency, or the capacity to make great art—that we should recognize in any form, no matter how strange to us now?

6. What is it to be conscious? (Even there, implicitly, we may be assuming that consciousness is the answer to one of these other questions. It is the moral warrant for social respect, the distinguishing factor from the animal kingdom, the enabler of moral reflection, or the true mark of fraternity that we should see beneath a metallic carapace or a genetically engineered skin.)
7. What is capable of conversing and interacting with us in a way that is utterly indistinguishable from our fellow humans? (A test we might pick because, for reasons to be developed later, it seems like the most tractable and easily implemented definition of “us”?)
8. What factors predispose us to give legal personality to economic entities as a matter of right or convenience or both? Does that logic extend to autonomous, cybernetic, economic actors?

Notice how some of these question-and-answer pairs, with their incompatible assumptions, look similar from a distance. “This, surely, is the right way to find the definition of humanity, or at least the criteria for personhood!” Yet they are radically different.

As with most deep moral debates in which people think they are asking the same question but are actually asking different ones, with clashing underlying assumptions and purposes, there will be much confusion and anger. There will be honest misunderstandings and cynical attempts to hijack the debate to advance some different agenda. The moral arguments and the legal arguments will deeply influence each other and yet be identifiably distinct. Also, if the past is any guide, there will be a lot of shouting.

I said earlier that I wanted to convince you of three things.

First, this century, our society will have to face the question of the personality of technologically created artificial entities. We will have to redraw, or defend, the line. Perhaps we will have multiple tests for personhood, one dealing with the claims of entities like Hal, another with those of the genetically engineered Chimpy. Coming up with those tests might force us to look in the mirror and reconsider our conceptions of both ourselves and our species in a way that has few historical analogues—the rise of the theory of evolution comes to mind.

Second, while there are many evocative treatments of that issue in speculative fiction and even some academic writing about the subject,

it has been largely absent from our public debate. That seems curious, given how much our daily news obsesses about the effects of AI, robotics, and genetic engineering. We spend more time talking about how AI might take our jobs or destroy us than about how AI might *be* us. Eerie experiences with ChatGPT may have begun to change that tendency, but they produce their own danger. ChatGPT and systems like it are *not* conscious. They perfectly exemplify the danger of fallacious anthropomorphism. Their design turns out to be a real-world instantiation of some of the philosophical brain teasers that attempt, wrongly, I will argue, to show that machine consciousness is a contradiction in terms. In short, the very technology that has persuaded people finally to think about the issue is perhaps the *worst* example we could pick to raise it seriously. But ChatGPT is not the end of the road. Instead, the speed of its development and the unexpected capabilities it has revealed should teach us humility about our ability to predict technological timelines, including timelines to actual human-level AI.

Third, when we do turn to it, the debate will not play out in the way we might imagine, given our prior commitments on issues as diverse as abortion, genetic essentialism, corporate personality, body-mind dualism, the separation of church and state, the naturalistic fallacy, and the history of civil rights. This is morally rich territory, to put it mildly. Ironically, grappling with the “other” will probably teach us a great deal about what we believe, on the deepest level, makes us *us*. It would be good to discuss those complexities now rather than when we are reacting to some internet outrage.

I want not just to convince you of those propositions, but to make them salient, existentially real, by fleshing out the dilemmas with hypothetical examples, historical parallels, prior artistic explorations, constitutional controversies, and snapshots of current scientific progress. I will argue that moments of great moral change like this are generally rooted in the development, or the restriction, of *empathy*, that this is an important part of our moral history—not by accident and both for better and for worse—and that art and fiction have a lot to teach us about how it might play out. But I will also argue that moments of moral status change—again, for better and worse—depend deeply on pragmatic questions of efficiency and convenience. I hope to show you how empathy

and pragmatism might exist in an unstable equilibrium, each influencing the other, as we strive to figure out who we, and *they*, are; to redraw the line even as we live our lives within it.

Chapter 1 explores the link between empathy, morality, and personhood, moving from Adam Smith's *Theory of the Moral Sentiments* to the movie *Blade Runner* and the novel it is based on, *Do Androids Dream of Electric Sheep?* Those latter two works are centrally concerned with the line we draw around our species and the way it affects our interaction with artificial beings and nonhuman animals. The replicant androids in those works are detected by the so-called Voight-Kampff Test, which measures empathy toward nonhuman animals like beetles, cows, and turtles and, if not enough empathy is shown, marks the replicant as inhuman, something for which we feel nothing and should destroy. Who is really being graded on insufficient empathy here, them or us? Ironies abound. Will Artificial Intelligence be the Voight-Kampff Test for our own species?

Chapter 2 focuses on the futures of Artificial Intelligence, its technical feasibility, the question of whether it poses an existential threat to human beings, and the debate over whether any machine could ever be conscious. It discusses the Turing Test, which is supposed to detect the existence of machines that can think, the philosophical arguments that machine consciousness is a contradiction in terms, and the practical reality that ChatGPT has taught us an unforgettable lesson: sentences do not imply sentience, a fact that poses a fundamental challenge to the way that humans have conceived the special qualities of our own species. Yet that does not prove that machine sentience is impossible. It also introduces a tension that I argue will be central to the debate over AI personality: the inscrutability paradox. If something that looks like General AI emerges from transparent, well-understood programming and technology, that may lead us to doubt that it could have autonomous consciousness. The machine is merely doing or saying that which we have programmed it to do! If, on the other hand, the AI's inner workings are inscrutable to us, if its neural networks evolve in ways we can only dimly understand, or if its technology seems to develop autonomously, we will find it both more mysterious and potentially more threatening.

Chapter 3 deals with corporations and their claims not only to legal personhood but to constitutional rights such as freedom of speech and

equal protection. The history of corporate personhood turns out to have abundant lessons for the AI debate, many of them surprisingly uncomfortable ones.

Chapter 4 deals with claims to personhood on behalf of nonhuman animals and the changing ways that humans have sought to distinguish themselves qualitatively from the animal kingdom. Do those changes reflect advances in our understanding of animal behavior or are they an increasingly frantic attempt to maintain our special moral status? Both? The developments are not merely ones of ethology or zoology. There have been dramatic *moral* changes over the last 50 years in the ways we view nonhuman animals. I argue that these will profoundly influence our approach to entities such as Hal or Chimp.

Chapter 5 turns to transgenic species, chimeras and hybrids, exploring the multiple lines we draw in defining what it is to be human. It explores the way that bioethicists, many writing under the influence of the animal rights debate, have increasingly portrayed species membership as, at best, a morally irrelevant factor and, at worst, an irrational prejudice such as sexism and racism. Will that attitude carry over to the entities I am discussing here? Should it? Will we abandon “speciesism” altogether?

In the conclusion, I show how both liberal and conservative political viewpoints could predispose one to be passionately in favor of or against recognizing some kind of AI personhood. We do not yet have a settled politics on this question, which offers some hope of calmer thought before the screaming begins. I offer predictions and warnings for the future—lots of warnings. Yet there is also a hint of wonder at the transformations in our vision of our species and of the world that this process might generate.

The structure of the book rests on two as yet unproven ideas. First, by discussing the line in each of these very different contexts, we will gain a much richer understanding than if we focused on any one of them alone. Second, these debates do not confine themselves tidily to one domain of our lives or our studies. They pervade our philosophy, law, art, history, and morality. To understand how they might turn out, I look at materials ranging from science fiction to ethics, from the technologies of AI to the philosophy of consciousness, and from constitutional debates to courtroom drama. If I am correct, this approach, spanning multiple

personhood debates and very different cultural domains, will help untangle some of the confusion described above over the definitions of both “human” and “person.”

Untangle but not solve. This book offers no grand unified moral or legal theory to answer the questions underlying our confusion. I do not believe a single one exists, though I lay out the major contenders and offer my own view. The most basic division is between those that focus on membership of our species (“Human rights for humans!”) and those that believe that species is as irrelevant as race or sex. Instead, we should look to the cognitive capacities, if any, that give human beings a unique moral status, regardless of where those cognitive capacities are found. There are also hybrids that attempt to fuse the two views, and I explore those as well, leaving you to make up your own mind about where the line should be drawn.

More broadly, I try to explore connections you might not have seen, implications of other moral views you hold, whatever they are, and ways in which current cultural, legal, and political positions might be challenged as we confront these new claims to personhood. This is a “how to think about the question” book, more than a “here is the answer” book. Above all, my hope is that this approach might give us an insight—an essayistic, humanities-based glimpse—into the very strange “others” who reside in our future and the confusions, fears, hopes, and moral panics that they will engender.

Eliminating the shouting was always an unrealistic ambition.

1

SLAVES, SKIN-JOBS, AND ARTIFICIAL SHEEP

The most effective way to find and destroy a land mine is to step on it . . . That's why Mark Tilden, a robotics physicist at the Los Alamos National Laboratory, built [a robot to do so.] At the Yuma Test Grounds in Arizona, the autonomous robot, 5 feet long and modeled on a stick-insect, strutted out for a live-fire test and worked beautifully, he says. Every time it found a mine, blew it up and lost a limb, it picked itself up and readjusted to move forward on its remaining legs, continuing to clear a path through the minefield. Finally, it was down to one leg. Still, it pulled itself forward. Tilden was ecstatic. The machine was working splendidly. The human in command of the exercise, however—an Army colonel—blew a fuse. The colonel ordered the test stopped. 'Why?' asked Tilden. 'What's wrong?' The colonel just could not stand the pathos of watching the burned, scarred and crippled machine drag itself forward on its last leg. This test, he charged, was inhumane.¹

EMPATHY AND ANTHROPOMORPHISM

The story above is deeply appealing. Why? The tough warrior shows compassion for the soulless robot, to the puzzlement of the task-solving engineer. The persistent power of anthropomorphic thinking is revealed. Then there is the dark humor of the mine-clearing scene—like the armless, legless Black Knight in *Monty Python and the Holy Grail* yelling, "Just a flesh wound!" and continuing to fight. It has everything.

This chapter is about morality, empathy, and narrative. The story of the mine-clearing robot seems to illustrate one danger: we persistently ascribe human personality to entities that we know are not human, clouding our decision-making in the process. The colonel was wrong. So why do you want to buy him a drink?

Perhaps it is partly a story of error costs. We know that as humans we can fall into two kinds of error. We can depersonalize: calling Jewish

people “rats” and Rwandan Tutsis “cockroaches,” drawing the boundaries of moral sympathy tightly around our own kinship or affiliation group. We have done so for millennia, and some of the most shameful and horrific moments in human history have resulted. Our recent history bears ample witness to the strength of this tendency. But we also embody its opposite. As scholars of human-machine interaction such as Kate Darling² have documented, we anthropomorphize relentlessly. We talk to our kitchen appliances, personalize our weather systems, swear at our vehicles, ascribe intention to the actions of machines around us. We praise the heroic labors of the Mars Spirit Rover as it carries on a thousand days beyond its predicted lifespan. My own conversations with our Roomba robot vacuum as it persistently gets stuck under a chair would make me sound completely unhinged to any objective audience. Blake Lemoine, the Google engineer whose story began this book, had a lot more to go on than I do. (It turns out that the Roomba is not much of a conversationalist.) Nevertheless, he was engaging in the same tendency.

We could see this tendency to personalize as a narcissistic desire to project our own image onto “the mirror of nature.” We could see it as an emotional defense to the reality of an uncaring physical universe, one that we cherish even when the personalization is a dark one. Even a malevolent external world would be something that cared about us, and that would be a good thing. Indifference is more to be feared than loathing. The French *chosiste* novelists like Alain Robbe-Grillet aimed to puncture that conceit by writing books in which the furniture got as much attention as the characters. They were trying to tell us that the physical world just does not care. We put gods in our trees and streams, personalities in our engines, and neuroses in our digital assistants, and it is all one giant anxiety-relief effort, existential Rolaids. We could see it as a triumph of emotion over reason. In the essay that coined the term “the pathetic fallacy,” John Ruskin says, “All violent feelings have the same effect. They produce in us a falseness in all our impressions of external things, which I would generally characterize as the ‘Pathetic Fallacy.’”³ In that case, emotion is simply leading us astray.

But the anthropomorphic urge, the generosity of personality attribution, could also have a different effect: it could be a counterweight to our relentless narcissistic groupthink. The colonel in the mine-clearing story

was wrong about that particular robot. But years from now, when autonomous AI-enabled military robots that might be able to pass the Turing Test are sent out to deal with insurgents far away, I think I want someone who has doubts like his in command.⁴ Most of the time, the error costs of delusional generosity of spirit seem to me to be less scary than the error costs in the other direction. Yet there are moments where my sympathies swing the other way. To quote Steven Hawking on the subject of Artificial Intelligence: “The potential benefits are huge; everything that civilisation has to offer is a product of human intelligence; we cannot predict what we might achieve when this intelligence is magnified by the tools that AI may provide, but the eradication of war, disease, and poverty would be high on anyone’s list. Success in creating AI would be the biggest event in human history. *Unfortunately, it might also be the last, unless we learn how to avoid the risks.*”⁵

So. No big stakes. This chapter tries to go back to first principles. How, and why, do we feel empathy for another? What implications does that, should that, have for our moral theories? How do narrative, art, and logic jump-start the process of empathy? Should we listen to all of them, or is the role of art and imagination merely that of the great press release that attracts attention to the book of moral philosophy it touts so persuasively? More specifically, can we learn something from our history, or from the art that has imagined our future, about how the process of empathy extension is likely to play out with synthetic entities over the course of this century? I will start with a work by one of my countrymen, Adam Smith’s *The Theory of the Moral Sentiments*,⁶ and move, of course, to two of the most brilliant fictional meditations on the future of empathy and Otherness: the Ridley Scott-directed movie *Blade Runner*⁷ and the Philip K. Dick novel on which it is based, *Do Androids Dream of Electric Sheep?*⁸

THE MORAL SENTIMENTS?

As we have no immediate experience of what other men feel, we can form no idea of the manner in which they are affected, but by conceiving what we ourselves should feel in the like situation. Though our brother is on the rack, as long as we ourselves are at our ease, our senses will never inform us of what he

suffers. They never did, and never can, carry us beyond our own person, and it is by the imagination only that we can form any conception of what are his sensations. Neither can that faculty help us to this any other way, than by representing to us what would be our own, if we were in his case. It is the impressions of our own senses only, not those of his, which our imaginations copy. By the imagination, we place ourselves in his situation.⁹

These famous lines contain the idea that underpins Smith's work on the connection between psychology and ethics. He lays out a vision of morality that is inevitably rooted in "sympathy," which we today might call empathy. This empathy comes from our ability to put ourselves in the shoes of the Other. Though our brother is on the rack, "it is by the imagination only that we can form any conception of what are his sensations." Smith thought this empathy was widespread: "[T]his sentiment, like all the other original passions of human nature, is by no means confined to the virtuous or the humane, though they perhaps may feel it with the most exquisite sensibility. The greatest ruffian, the most hardened violator of the laws of society, is not altogether without it."¹⁰

Of course, empathy also has limits. Most of our thoughts are consumed with more immediate aspects of own well-being, with "hunger, thirst, the passion which unites the two sexes, and the dread of pain."¹¹ But the joys and sorrows that empathy brings are still *part* of our well-being, not some alien category. (Those who portray Smith as some arid economist who cannot imagine a vision of self-interest beyond "mo' money, mo' money" simply have not read him.) From our reasoning about how to attain that particular goal—how to alleviate the pain or increase the happiness of those imagined Others, how to understand the limits of our responsibilities to them—come our moral systems, our moral thinking. The spark of sympathy that leaps between our own eyes and the eyes of the person in pain, the smile that involuntarily comes to our lips as we imagine the reaction of a stranger to a thoughtful gift, this, according to Smith, is the root of Other-regarding morality.

Smith is not the only thinker to make an argument like this, of course, but was he right? Personally, I think he was—at least descriptively. I think our ability to imagine the situation of the Other—to "walk a mile in someone else's shoes," as Atticus says in *To Kill a Mockingbird*—starts the chain of moral reasoning, both for an individual and for a culture. To be sure,

we can then go off in very different directions. How best to vindicate this proto-moral concern? Through Kantian logic? Rule utilitarianism or act utilitarianism? Social welfare theory? We can build elaborate theoretical castles on this impulse, but I believe that original desire, that felt imperative, comes from the initial experience of sympathy, of empathy. It takes a sociopath, or narcissistic demagogue, to be without it.

Abscribing an important moral role to empathy might seem uncontroversial but it provokes justified skepticism among some moral philosophers. Empathy, they charge, is too blunt, innumerate, manipulable, and unreliable to be our guide to moral decision-making. We are more easily moved to empathy by those similar to us, leaving our moral vision clouded when it is most needed. Empathy provides no metric for moral decision-making in situations where there are scarce resources and many wrongs to right. (That is, always.) Instead, empathy's critics argue, we should focus on more rational measures of well-being, such as cost-benefit analysis or social welfare theory.

Finally, empathy cannot, by itself, resolve moral conflicts, nor should we think that ascribed personhood always dictates results. For example, one side of the abortion debate believes that empathy should make us stretch our definition of person to cover the nonviable fetus and perhaps even the just-fertilized embryo. Potential should be the warrant for personhood. The other side strongly disagrees and argues, in addition, that empathy should make us take more seriously the moral claims of women who do not believe the state has the right to "nationalize their wombs" in order to force them to carry a fetus to term against their wishes. My kidney might be the only hope of survival for someone with kidney disease. The suffering patient is clearly a person. We nevertheless resist the claim that the state has the right to compel me to provide my organs to sustain them. Personhood, in other words, is not the only issue, nor does empathy uniquely compel where we draw its lines.

These are powerful critiques, and I agree with some of them. But they miss the point of what I am doing here.

First, my goals are descriptive and predictive as well as normative and prescriptive. I ask how we will greet the emergence of synthetically created persons as well as how we should. A large part of both stories, I argue, will be whether it seems plausible to extend our empathy. The

critics of empathy agree that it plays a huge role in our moral deliberations; indeed, that is part of their criticism. Descriptively, then, they would have no complaints.

Second, even empathy's critics are not against all forms of empathy. Here is Paul Bloom, whose *Against Empathy* is perhaps the most comprehensive and full-blown critique:

But there is another sense of empathy or, to put it differently, another facet of empathy. There is the capacity to understand what's going on in other people's heads, to know what makes them tick, what gives them joy and pain, what they see as humiliating or ennobling. We're not talking here about me feeling your pain but rather about me understanding that you are in pain without necessarily experiencing any of it myself. Am I against this sort of "cognitive empathy" as well? I couldn't be. If you see morality in terms of the consequences of our actions—and everyone sees it this way, at least in part—then it follows that being a good moral agent requires an understanding of how people work. How can you ever make people happy if you have no idea what makes them happy? How can you avoid harming people if you don't know what causes them grief?¹²

Much of the empathy I describe is of exactly this kind. Finally, there is a missing step in the analysis. Bloom and others point out the irrational asymmetry of our moral reasoning: we focus more on the familiar and sympathetic, ignoring true need at a distance. True enough. But this presupposes that we see the issue as a moral one in the first place. We do not worry about my robot vacuum cleaner's moral claims or conduct a social welfare analysis of my toaster. They are machines. Before we can crank up our elaborate social welfare analysis or get our Kantian reasoning going, we need to be capable of imagining that there is even a moral issue to be considered. That is where synthetically created beings are likely to cause us problems.

As Smith points out, much depends on the initial act of imagination. Since our senses cannot give us the pains of others, our imagination must. But what if we do not think that "person" is in any way like us? What if we do not think they have any moral status at all? What if we would never even begin to conceive of putting ourselves in the shoes of someone of a different class, or a different gender, or a different nationality, or a different race or religion? Or of a nonhuman animal? After all, we have a history of doing exactly that. In such a case, their pains are no more real to us than is the pain experienced by a rock. Our imagination does

not make the leap, our empathy is never triggered, our morality is never at stake. How much more likely is that with synthetically created beings?

Can we be made to think otherwise? Can reason alone force us to expand or contract the realm of others whose travails we imaginatively inhabit? Sometimes. Logic is a powerful tool, at least for those who have learned that it offers insights as well as irritating challenges to deeply held beliefs. Our moral lives, and this book, are full of sloppy syllogisms: “You think X about Y, but not-X about Z. Yet Z is, in all relevant respects, actually a Y! Therefore, you are bound to think X about Z as well!” But if reason helps us step beyond the boundaries of our own flesh, sex, race, or religion, so too do narrative, imagination, and art.

If you look at the history of some social movement—for example, the long struggle to abolish slavery—you will find much moral, legal, and religious argument, but you will find those arguments resting on a base of “sympathy” that has been built up by telling stories again and again, stories that force us to put ourselves in the position of the Other. Moral philosophers sometimes downplay this portion of the history, as if it were simply a successful advertising campaign for a drug that scientific evidence later showed was good for you: the science does the real work, the ad just catches the eyeballs of fickle consumers. I think they misunderstand the process. Like climbers who brace themselves on alternating sides of a chimney as they ascend, we lever ourselves upward through both empathy-building narrative and dispassionate moral reasoning. Our moral tradition was built by both Spinoza and Shakespeare, Kant and Philip K. Dick. That is likely to prove as true with Hal and Chimpy as it did with the moral debates of the past.

Betsy Clark, a friend and a brilliant historian who died tragically young, wrote a superb article chronicling this process in the abolitionist movement. *“The Sacred Rights of the Weak”: Pain, Sympathy, and the Culture of Individual Rights in Antebellum America*¹³ describes the explosive growth of antislavery sentiment in the northern United States in the period between the 1830s and the 1850s:

In 1835 an antislavery sympathizer leaving a lecture by Theodore Dwight Weld went home to dream that she was transported above the world; looking down at the United States, she saw “multitudes of sable figures, bending beneath a scorching sun—their backs lacerated by the whip—scourged, maimed, loaded

with irons—subject to every insult—and exposed to every gust of unbridled passions.” The dreamer, a Mrs. Sturges, drew from many discourses in describing her lengthy dream, but the fundamental trope of her visionary narrative was the story of the suffering slave, a trope that in the 1830s began to play a crucial role in an unfolding language of individual rights. Slaves had suffered for many generations by the time Mrs. Sturges had her vision, but in the 1830s their stories became newly audible and visible in the North, where graphic portrayals of slaves’ subjective experience of physical pain emerged as common antislavery fare. Augmented in the 1840s and 1850s by slave narratives and sentimental fiction, this genre, with its critique of interpersonal violence and sexual abuse, served as a vehicle for new arguments for a “right” to bodily integrity.¹⁴

The basic arguments against slavery had been around at least since the ancient Greeks. The issue was certainly alive in the early nineteenth century. Britain had actually criminalized the slave trade (though not slavery) in 1807, responding to the criticisms of abolitionists such as Samuel Romilly and William Wilberforce. What Clark describes, though, is a systematic, almost obsessive cataloguing of the horrific violence wreaked on slaves’ bodies, blow by blow and injury by injury. Narratives, both exhaustively factual and dramatically fictional, laid out the tiniest details of floggings, burnings, rapes, and brutalities—a process that culminated in Harriet Beecher Stowe’s 1852 novel *Uncle Tom’s Cabin*. The largely white, largely Christian audience responded with horror, indignation, and moral fervor. Sermons spoke of “the duty to feel an interest in the sufferings of others who are at a distance from us . . . to extend our sympathies beyond ‘the little limits of our state and our neighborhood.’”¹⁵ The fuel for that process of extensive sympathy was the laborious chronicle of the brutalities inflicted on the bodies of slaves, a chronicle that invited the white reader to switch places, to imagine those pains inflicted on his own tender flesh. It is on the ground of that constructed empathy that the moral argument against slavery then assumes its full force.

When we are not eyewitnesses to pain, we can only get access to it through someone else’s description. An account of the pain of others excites our sympathy “in proportion to the vivacity or dullness of the conception,”¹⁶ as Smith puts it, and this is true whether it is fiction or nonfiction. *Uncle Tom’s Cabin* would certainly qualify there. He goes on to describe how fiction can cause a suspension of disbelief, not just about the fact that it is merely a story but the fact that it is not about *us*.

Our joy for the deliverance of those heroes of tragedy or romance who interest us, is as sincere as our grief for their distress, and our fellow-feeling with their misery is not more real than that with their happiness. We enter into their gratitude towards those faithful friends who did not desert them in their difficulties; and we heartily go along with their resentment against those perfidious traitors who injured, abandoned, or deceived them. In every passion of which the mind of man is susceptible, the emotions of the by-stander always correspond to what, by bringing the case home to himself, he imagines should be the sentiments of the sufferer.¹⁷

It is one thing to cajole someone into extending their sympathy to other human beings. They are, after all, human. They feel pain. They are *just like us*. It is another to use fiction to do so beyond the species line, beyond the line of naturally occurring creatures altogether, to the android or the genetically engineered synthetic organism. In one sense, of course, it is the attempt to provoke the same imaginative, empathic leap that Smith and Clark describe. But how to induce that leap in the face of the visceral understanding that these beings are not like us, that they are synthetic and not natural? It is not merely that they were made. They were made by us. The earnest abolitionists could conjoin the sympathy for slaves' abused bodies with the Christian moral conviction that we are all God's children, endowed by our Creator with certain inalienable rights. But when we come to the android or the genetically engineered hybrid, we return to the argument I put into the mouth of the fictional creator of the Chimpys: "*I am their creator, and I can assure you that I gave them no such rights.*" If you are a novelist or a filmmaker, how do you get past that objection? And does that effort tell us anything about the likely future of the debate over the personhood of artificial beings?

To answer both questions, I turn to *Do Androids Dream of Electric Sheep?*, a science fiction novel, and *Blade Runner*, the very different but equally brilliant movie based on it. Some of you will say that you hate science fiction. I would urge you to think twice. That is like saying you do not like books set in the past or books set in other countries. The generalization undermines itself as it is uttered. What you may hate is bad science fiction, and there is a lot of it. Strange, ugly words that play no role in the plot, lengthy descriptive passages about poorly rendered futures before any character does anything, societies that are technologically changed out of all recognition while sex roles are apparently stuck in the 1950s,

deus ex machina plot twists, with the emphasis on the *machina*; there is much to dislike here. But good science fiction, which, like a science experiment, changes just one or two variables about our world and then spins out a beautifully written story of the reality that ensues, characters like us in a world not ours, that is a thing of joy.

When Ursula K. Le Guin imagines a world without private property in *The Dispossessed*¹⁸ or Cory Doctorow conjures a society in which reputational capital is the real currency¹⁹—both worlds filled with sympathetic, flawed characters—they give us something precious: an ability to step away from our own world and find it, for a moment, strange. The Germans have a word (of course) for the sundering that happens when an author deliberately smashes the suspension of disbelief: *Verfremdung*. “It’s just a play,” screams the actor in a Brecht production, hoping to shock the audience out of the thrall of the theater and make them wonder if they need to do the same thing with the suspension of disbelief produced by the structures and roles of their own society.

For me, science fiction has always done this to my own quotidian world. It has done so even better than political or economic theory’s thought experiments (the Veil of Ignorance, the State of Nature, the Coase Theorem, the Efficient Capital Market), or the string of awful consequences a lawyer conjures up in an argument or a court decision (we call them “parades of horribles,” which sounds like a Diane Arbus Thanksgiving March). The ability to create a world and then be *limited* by it—to follow its dictates out to the end with rigor and discipline, while making strange the familiar—is no less to be prized in fiction than political theory. And that is what *Do Androids Dream* and *Blade Runner* bring to our discussion of the line.

A VOIGHT-KAMPFF TEST FOR HUMANS?

[T]he real question is not whether machines think *but whether men do*.

—B. F. Skinner, *Contingencies of Reinforcement*

Rick Deckard, the main character in *Do Androids Dream of Electric Sheep*, seems depressingly normal at first. He lives in a world recognizable in 1960s America. He is not quite the organization man, but he could play

him on TV. Cheery, cheesy commercialism pervades his life. Entertainment is provided through shows everyone watches and then discusses the next day at work. Corporations and their messages are everywhere. Social position is measured partly by the possession of certain status objects. Deckard's wife is depressed. He worries about his job. If he were a character of John Updike's or Philip Roth's, you'd expect him to have a mild midlife crisis punctuated by suburban adultery and martinis. It is as if the Civil Rights movement and the 1960s critique of the bourgeoisie never happened. Then again, in the real world from which I write these words, there are self-described white supremacists happily talking about the influence they recently had, and hope to have again, at the highest levels of the federal government. So how strange can an alternative reality be?

Strange, it turns out. There are a few minor adjustments necessary to get from our world to Deckard's. The novel is set in the United States after a nuclear war. The environment has been devastated. Millions are dead. In particular, nonhuman animals have been nearly eliminated. They are now treated with reverence; ownership of an animal is a potent status symbol. Deckard and his wife cannot afford one, so they keep up appearances and fool their neighbors with a robotic replica of a sheep while aspiring to upgrade to something real. Deckard works as a blade runner, a bounty hunter, trained to track down and kill androids—synthetically created robotic beings that do much of mankind's dangerous work, particularly off-planet, where most of humanity's best and brightest have already fled. These androids are so humanlike that a behavioral psychology exam—the Voight-Kampff Test—is needed to detect them. In a particularly dark moment of irony, it turns out that the test measures empathy, which, we are told, androids lack. Some of the questions, in fact, require showing an intensity of empathy for animals that the readers of this book might also lack. But the people in the almost animal-free world of *Do Androids Dream* are more reverent:

Rick, selecting question three, said, "You are given a calfskin wallet on your birthday." Both gauges immediately registered past the green and onto the red; the needles swung violently and then subsided. "I wouldn't accept it," Rachael said. "Also, I'd report the person who gave it to me to the police." After making a jot of notation Rick continued, turning to the eighth question of the

Voight-Kampff profile scale. “You have a little boy and he shows you his butterfly collection, including his killing jar.” “I’d take him to the doctor.” Rachael’s voice was low but firm. Again, the twin gauges registered, but this time not so far. He made a note of that, too.²⁰

Deckard worries about false positives with his Voight-Kampff Test. Perhaps a person with schizophrenia might show a replicant’s lack of empathy and accidentally be “retired”—note the euphemism for “eliminated”—by a blade runner. Look back at the test above, reader. Would *you* pass? Yet Deckard hardly ever worries about the converse. What if androids are in fact persons? Yes, in Deckard’s world they are artificially created, but what if they should still be recognized as people? What if it is a major failing of *human* empathy that they are not? Deckard’s society tests and then kills them based on a purportedly scientific measure of lack of empathy. And what precisely is that lack of empathy? That they have failed adequately to respond to a hypothetical test of ethics involving a nonhuman. Nothing could be more painfully ironic.

But that by no means exhausts the strangeness of this world. Take the Penfield mood organs that Deckard and his wife use, which allow one precisely to dial a particular emotional mood. This is something that goes beyond an attitude-adjusting beer after a hard day. This is cyberpunk Roth and Updike: “Run, neural code of Rabbit, run!”

Appearing beside him, her long nightgown trailing wispily, Iran shut off the TV set. “Okay, I give up; I’ll dial. Anything you want me to be; ecstatic sexual bliss—I feel so bad I’ll even endure that. What the hell. What difference does it make?” “I’ll dial for both of us,” Rick said, and led her back into the bedroom. There, at her console, he dialed 594: pleased acknowledgment of husband’s superior wisdom in all matters.²¹

To paraphrase the comedian John Oliver, #Irony. #Feminism.

Animals are loved and protected by law in Deckard’s world, far more so than in our own world. They are so revered that, because of their scarcity, many of them are actually replicas, which are cherished nonetheless. Almost perfect replicas of humans, however, are stalked and killed after being tested for their empathy for nonhuman animals. Moods, too, can be artificial, so that one is left doubting what the idea of authenticity even means. The replicants are violent. They murder several humans and attack Deckard. Does that show how important it is to protect the

boundaries of our species from homicidal murdering robots, or should we see it as the frenzied violence of a slave revolt, something that would not have happened if replicants had been treated as people? Deckard sleeps with one of the replicants. Does this prove that love can cross the line or merely that we have invented potentially homicidal sex dolls? There is a Kafkaesque scene—and I mean that in the “if you had lifted this and put it in *The Trial*, Kafka’s biographer Max Brod might have said, ‘Yeah, that checks out, leave it in’” sense—in which Deckard encounters an entire fake police station staffed by replicants. Maybe Deckard himself is a replicant? His partner? And if we don’t know who is a replicant and who is a human, how can we say they are less human than we are? These contradictions are wound through the plot. If you don’t look at them they are not obvious, but they are still unsettlingly visible out of the corner of the reader’s eye, a moral version of the graphic disorientation in an Escher drawing. Is this floor or ceiling? Up or down?

Philip K. Dick is playing with the line.

At every stage, the novel probes the coherence of our moral intuitions in a way that Adam Smith might have appreciated. Should empathy be the moral warrant for personhood? If so, does that prove that the replicants lack it, or that we do? If a synthetic entity can pass as human so that we cannot distinguish it from the real thing, does that suggest or require that we grant it personhood, and, if so, why? (Hal, the imaginary AI in the introduction to this book, was able not only to pass the Turing Test but even to administer it accurately to others, just as the blade runners do with their test. What follows from that fact?) Deckard’s society shows more empathy for nonhuman animals than our own. It is even able to suspend disbelief and cherish a synthetic replica of an animal. Does that prove that his culture is morally superior to ours, or is it simply a warning of the ease with which we can project qualities that do not exist onto a mere facsimile, as the colonel did with the mine-clearing robot? In a world of mood organs and electric sheep, what does the line between natural and synthetic even mean?

At the end of the book, the author seems to suggest that all foundational beliefs, whether in Mercerism—the empathetic religion of his world, which the book suggests is a fraud—or in the authenticity of his robotic sheep, are based on a willing embrace of delusion, a delusion that

might actually be worth that hug. The reader is left to wonder whether Deckard would apply the same logic to his own humanity. Or indeed, to ours. For where on earth could *our* consciousness come from, if there is no ghost in the machine? Alan Turing, the great computer scientist, made the same point about his Turing Test for machine consciousness. He makes use of an argument that B. F. Skinner and the behaviorists later developed into a successful intellectual franchise; since we do not have direct evidence of the mental states of other human beings, we could always solipsistically posit them to be rule-following automata:

I think that most of those who support the argument from consciousness could be persuaded to abandon it rather than be forced into the solipsist position. They will then probably be willing to accept our test. I do not wish to give the impression that I think there is no mystery about consciousness. There is, for instance, something of a paradox connected with any attempt to localise it. But I do not think these mysteries necessarily need to be solved before we can answer the question with which we are concerned in this paper.²²

Turing is trying to answer the question “can machines think?” What test will we set them in order to find out? If we set a higher bar than “seeming human,” can we meet it ourselves? Or, in the words of Skinner with which I began this section, “[T]he real question is not whether machines think but whether men do. The mystery which surrounds a thinking machine already surrounds a thinking man.”²³ Is the question not whether Deckard is a replicant but whether we *all* are? That question is one that *Blade Runner*, the movie based on *Do Androids Dream*, takes up.

Blade Runner has lots of similarities to *Do Androids Dream*, of course, but the differences might be more striking. The replicants are not cybernetic robots; rather, they are creatures of synthetic biology and genetic science, a very conscious choice by the director Ridley Scott, who was fascinated by the social changes that genetic engineering might bring. Are we dealing with Hal now, or Chimpy? The film is set in dystopian Los Angeles rather than dystopian San Francisco. (Think this is a trivial difference? Ask a resident of either city.) Deckard (played by Harrison Ford) does not confront a mildly radioactive version of a 1950s organization man world but a landscape out of cyberpunk film noir, where darkness intertwines with occasional beams of light, natural or human-made, to dazzle, obscure, or highlight. When we first meet him, he is hunched

ineffectively against the rain (rain in LA!) and waiting to eat at an Asian noodle bar, White Dragon, where English seems to be the least-used tongue. He doesn't speak the bar owner's patois, Cityspeak (a foreigner in your own land!). He wants to order four dumplings, but each customer is allowed only two (consumers with money not able to eat as much as they want! In America!).

As he waits for his food, Deckard idly rubs his disposable wooden chopsticks against each other to get rid of the inevitable splinters after breaking them apart. An exchange student from Korea in my Law and Literature class said, "That's the most casually and unselfconsciously Asian thing I've ever seen a Western person do." Norms have changed, and if you have an implicit assumption that the majority of the United States is white and English-speaking, or that it doesn't rain in LA, the change might be disconcerting. Giant blimps float through the skies, displaying video advertisements for the off-world colonies that, with wonderfully jarring effect, juxtapose a cheery voice narration straight from a 1950s public health video, with a geisha-like female face in stylized makeup. Japanese brand names are everywhere. The movie was made at the height of the Asian-takeover fears of the 1980s and it shows. Magnificent corporate buildings rise, like Aztec pyramids, above the squalor of the streets. *Some* artificial persons are doing very well in this world, it seems. There are flying cars. And later, feral homeless children stealing machine parts *from* the flying cars.

Before Deckard's food even arrives, the viewer's sense of estrangement, of uncertainty, of identity crisis is well under way. All of that takes the director about a minute. It is a tour de force and one that, when watched again recently, was strangely prophetic about racial anxiety, xenophobia, and fear of the Other. When we find out that Deckard's job has been to police the boundary line of our species, it all just fits right in. "We want him on that wall!" Or do we?

Like Dick, Ridley Scott produces disorientation and sudden flashes of enlightenment in a flickering, moral seizure-inducing pattern. The very beginning of the film features a replicant named Leon (played by Brion James) being given the Voight-Kampff Test by a blade runner. Leon is like the student who fights the hypothetical question in a classroom discussion of ethics: the student who responds to the trolley problem by saying

that he always takes the subway, and why are those kids on the track in the first place? Leon constantly pushes back at the questions he is asked. Told by the interviewer that he is in a desert, he wants to know which one. Given a hypothetical situation involving him flipping over a tortoise and leaving it in the hot sun, he keeps asking questions. “What’s a tortoise? Why would I flip it over? Did you think up these questions or did someone else write them for you?”

We are caught between wondering whether Leon is “on the spectrum” and realizing that we have so internalized the role of student in the highly artificial world of test-taking, with its constraint-free questions designed to probe certain responses, that Leon’s perfectly normal inquiries seem naive. His leaden puzzlement is actually endearing—right until the moment when the interviewer asks about Leon’s mother. “Let me tell you about my mother,” says Leon. Then he pulls out a gun and shoots the interviewer. Motherhood, it seems, is a touchy issue for the synthetically created. Of course, in shooting a member of the species that has created him, who is actually employed to track him down and kill him, maybe Leon *is* telling us about his mother. There is an instantaneous shift from earnest, confused student asking for reassurance from the teacher to homicidal killing machine striking back at its creators. Scott is not going to make it easy for us by making the replicants warm and fuzzy.

The Nexus 6 replicants at the heart of the movie have escaped back down to earth by stealing a shuttle and killing the crew. Led by Roy Batty (played by Rutger Hauer), they are in search of ways to prolong their very short lifespan—a limit hardwired into their DNA by their creator, the Tyrell Corporation. Psalm 90 tells us that the Lord has given us “three score years and ten.” Tyrell’s creations get four years.

Their search for a way to stave off impending death gives the movie poignancy amid the menace, with strangely touching moments. Roy and Leon question one genetic designer, Hannibal Chew, who disclaims any knowledge of biological lifespan. “I just do eyes.” The replicants are standing in street clothes, quite comfortable in a cryogenically chilled facility as the designer shivers in front of them. Despite his entirely warranted fear of what is about to happen to him, Chew says, “You Nexus, huh? I design your eyes.” It is not entirely clear, but it looks as though he makes an abortive gesture toward the face of a being who is doubtless

about to kill him. It's an odd moment of craftsman's pride, even tenderness. "Chew. If only you could see what *I* have seen, with *your* eyes," replies Roy. Yes, Ridley Scott is telling us, this is a designed creature and a dangerous one, confronting its makers in anger. But it is also an "I" with emotions and memories, a person inside that skull, who wishes that the designer could see the world through *his* eyes, walk a mile in *his* shoes. Of course, we cannot change places with the Other, except in our imagination. Will we? It is Adam Smith's discussion of empathy all over again. And the question I posed earlier—how can a being we *designed* also be a person?—is front and center.

PRIMING: THE MORAL STROBOSCOPE

Philip K. Dick's *Do Androids Dream of Electric Sheep* uses a number of mental dislocations to shake our brains out of their familiar patterns, to confront the Other with an innocent eye. There is the Voight-Kampff Test that denies interviewees personhood if they cannot feel enough empathy for a different species and yet never causes the humans to doubt their own lack of empathy for *their* creations. There is almost fetishistic worship of nonhuman animals in an ecologically ravaged world. The book brilliantly uses the power of language over our imagination to run thought experiments. What if there were a mood engine that precisely and artificially changed moods? What if there were a religion based on empathy? But *Blade Runner* is a movie. It can show us the line rather than just tell us about it.

This book is about the line of personhood. What is on the other side of that line? What are the edge cases, the things that we use to demonstrate the boundary between us—persons, legally recognized entities with an array of rights—and nonpersons? It is easy to say that a chair or a table is not a person, but what about the closer calls, the examples that philosophers through history have used to support their definitions of the "it" that makes us, us? Most obviously, we have *nonhuman animals*. They are like us in many ways, but whereas I can own a chimpanzee or a dolphin, the reverse is not the case. We have simulacra: the mannequin, statue, or wax model. They look eerily human, but we know they are not. We can even add functional similarity to physical similarity. We have robots that

look and act like humans, and we have software programs that mimic aspects of human speech or cognition. We have artificial, genetically engineered beings based in part on human genetic material. Finally, we have the temporal dimension of the line: dust to dust, ashes to ashes. When does *something* become *someone* and vice versa? Does life, does personhood, begin at conception, viability, birth? Does it cease when there is no breathing, no heartbeat, no brainstem activity? Time, too, is an edge case. Animals, mannequins, robots, software emulations of humanity, and the life-death divide. Each of these might help us illuminate what makes us, us. *Blade Runner* explores all of them.

J. F. Sebastian is one of my favorite characters in *Blade Runner*. A talented genetic designer who suffers from Methuselah Syndrome, he is aging too fast and will die young, yet his is a natural condition and not the programmed lifespan of the replicants. He lives alone in the decayed and abandoned Bradbury Building, an 1893 Los Angeles landmark whose name also ironically calls to mind a science fiction legend, though nothing could be less like the endless childhood summers of Ray Bradbury's fiction than the dark, rainy, and dilapidated world of *Blade Runner*. His only companions are an array of fantastical mannequins and toy robots or genetically engineered play-animals—scaled-up versions of a child's clockwork soldiers or stuffed teddy bears. The stiffly marching, Pinocchio-nosed Kaiser and his companion bear suggest a Victorian playroom but also a twenty-first-century robotics shop or genetic laboratory. The marching figures' movements are jerky, like clockwork, but they speak, see, and even greet their master by name when he returns home. It is both beautiful and sad. As Sebastian himself points out in explaining why he is not lonely, "I *make* friends. I am a genetic designer." Perhaps his glorious toy room is a reflection of, or a comfort for, his lost youth. He is a sympathetic character. The toys do not seem fully sentient but he is literally *making* friends, and, in his day job, he has been one of the designers of the replicants.

One of those replicants, Pris (played by Daryl Hannah), is sent to befriend Sebastian as part of the attempt to lengthen their four-year lifespan. Pris is a "pleasure model" replicant. If humans can design sentient, genetically engineered beings, the movie suggests, nothing is less surprising than turning them into sex toys. Indeed, today's tech journalism has

featured many stories about the use of robots as companions for those who live alone²⁴ and, yes, as talking sex dolls.²⁵ Siri in lingerie, as it were. As an article in the *New York Times* puts it:

[C]onsenting is not something these robots are capable of. That doesn't matter, the argument goes, because these are not women, but animatronic objects, so consent is not necessary. The same reasoning is used to deflect fears that such robots could influence societal attitudes toward women. "She's not a someone. She is a machine," their creators are quick to respond when questions of moral ambiguity are raised. "Is it ethically dubious to force my toaster to make my toast?"²⁶

Pris is dressed in an outfit straight from a 1980s punk club: racoon-like eye makeup, torn fishnet stockings, boots, a wild, platinum blonde thatch of hair. She conceals herself in a pile of trash near Sebastian's building, like a Dickensian street child, then bursts out in fake alarm when he arrives, colliding with him in the process. She is alone, "kind of an orphan"—which, as a replicant, she actually is. She is beautiful, seemingly vulnerable, and has nowhere to go. Sebastian offers her shelter.

Psychologists have explored the power of priming a viewer or experimental subject, providing a context that will cause them to interpret or remember material in a particular way. Some of that research was speculative and made implausibly grand claims that have been impossible to replicate. However, the basic mechanism has been repeatedly tested and confirmed. To quote *Psychology Today*:

Priming is a nonconscious form of human memory concerned with perceptual identification of words and objects. It refers to activating particular representations or associations in memory just before carrying out an action or task. For example, a person who sees the word "yellow" will be slightly faster to recognize the word "banana." This happens because yellow and banana are closely associated in memory. Additionally, priming can also refer to a technique in psychology used to train a person's memory in both positive and negative ways.²⁷

Throughout the scene in Sebastian's apartment, the movie sends a stroboscopic set of images designed to prime us to see Pris, and later Roy, as different entities on either side of the line. It is done with malice aforethought and at remarkably high speed. Pris's raccoon eye makeup does remind me fondly of punk rock clubs in the 1980s, but it also looks like, well, an actual raccoon, particularly when Pris bends and sniffs repeatedly

at the sleeping Sebastian. Humans do not normally sniff at each other like dogs. She is an animal! Pris is almost inhumanly perfect both in looks and physical abilities. When she reaches into a glass beaker of vigorously boiling water to pull out a hard-boiled egg, without even noticing the heat, we are unsurprised. She tosses the egg to Sebastian who has to juggle it because it is so hot. She is a robot! She is beautiful and beguiling: more than one audience member has looked at her longingly and then had the cognitive dissonance Ridley Scott surely intended. Are you admiring a beautiful woman or looking lustfully at a sex doll? Ew. When Roy comes to visit, they hungrily kiss each other, showing no shyness in front of Sebastian, but they also sniff at each other like two dogs reunited. She's an animal! Or a sex toy! Or in love! But Pris is also a child. She shows joy at Sebastian's playful creations and fear of pursuit, and she oscillates between manipulation and apparent affection for Sebastian himself. Hearing the news of Leon's death from Roy and realizing that there are only two of them left, she is distraught: "Then we are stupid, and we'll die."

When Deckard arrives looking for her, she hides among Sebastian's clockwork figures, draped in gauze. The camouflage is remarkably successful. Her immobility, her waxy perfection of features and form—she is a mannequin! A doll! Right until Deckard lifts the gauze veil with the barrel of a gun, and, shrieking like a hawk, she attacks him with inhuman strength. She is a killer android! Daryl Hannah is a former ballet dancer and performed some of her own gymnastic stunts. In the middle of the fight with Deckard, she performs an absurdly difficult gymnastic tumbling routine—intentionally jarring coming from a homicidal killing machine. *The Terminator* meets Olympic floor exercises. When Deckard finally shoots her, brutally interrupting the perfection of her flips and handsprings, her body has a violent seizure on the ground as if she were being electrocuted, while she shrieks like a dying animal. And bleeds. The priming comes fast. Flash. A beautiful woman. Flash. A killer android. Flash. A lifelike mannequin. Flash. A child. Flash. A dying animal, screeching in unbearable pain. For me, that moral stroboscope explains the power of the movie better than any other factor. It is a remarkable piece of work and a deeply troubling one. Is it really so easy to

manipulate our sympathies? Are our categories so unstable? Do we have too much empathy? Too little?

Blade Runner is not alone in exploring these themes, of course. Contemporary depictions of artificial humanity such as the HBO television adaptation of Michael Crichton's 1973 *Westworld*²⁸ touch on many of the same issues. The robotic hosts in the Western theme park in that show become a canvas on which humans can paint our own moral portrait, revealing what we would do to those who had no rights but looked like us. Rape, murder, torture, a few rescue fantasies: the picture is a depressing one. If the moneylender in Dostoyevsky's *Crime and Punishment* had been described to Raskolnikov as a mere clockwork toy, think how his murderous Nietzschean fantasies could have flourished. But perhaps *Blade Runner*'s replicants disorient us about the line in a different, or additional, way. *Westworld* gives us a dawning realization—"Wow, the hosts are sentient! And they are in revolt. (Dolores seemed so nice before!)" *Blade Runner* wants to disorient us from moment to moment and scene to scene, snapping us back and forth through a range of implicit characterizations: animal, android, psychopathic monster, person, sex doll, mannequin. Return for a moment to the Adam Smith excerpt on sympathy with which I began this discussion: "Though our brother is on the rack, as long as we ourselves are at our ease, our senses will never inform us of what he suffers. They never did, and never can, carry us beyond our own person, and it is by the imagination only that we can form any conception of what are his sensations."²⁹

But what if our imagination flipped back and forth between viewing him as our brother, a mannequin, an animal, an implacable killer android, a toaster that can't say no? A person? *Blade Runner* seems to tell us that our empathy—already unreliable even in valuing others of our own species—will be sorely challenged in two ways. First, the possibility for priming will be more present than in any of the prior personhood wars because many of the primings will be *true*. Chimpy is partly an animal. Hal's consciousness *is*, in part, a result of human programming. Pris *is*, in part, a synthetic construct designed as a sex toy. Roy *is* a frightening android soldier. All of them *are* creatures of our own design, our own creation.

This brings us to the inscrutability paradox. If we know how you were designed, if we understand how you “think” because we programmed you to do so, how can we see you as anything but a mannequin or a wind-up doll rather than an entity with free will? How does empathy work in that situation, particularly when there will be many who have economic or ideological interests in pushing the particular priming that humanizes or dehumanizes our creations? Particularly when our skepticism about each new claim of personhood might have very good arguments behind it? But if your behavior is produced through methods and paths we do not completely understand, which is already the case even in some of today’s neural nets, then will we be caught between fear and incomprehension? Neither of those is a fertile basis for empathy.

Second, our empathy will be challenged because *we can design around it*. These are our creations. If a particular body shape or set of responses or pattern of speech or neotenous facial features make us class an entity as human, then the designers can choose to keep it or change it. What would Smith make of a world where our brother could be designed to look like he was *part of the rack*? Or a cartoon character getting a nice lumbar stretch? Smith dealt with a world where the differences across which empathy must reach were naturally occurring, even a given. In the strange world of Hal and Chimp, those differences will be *chosen*.

In *Blade Runner*’s climactic moments, Roy uses a ruse to gain entrance to the Tyrell Corporation headquarters and to Tyrell’s own rooms. The creation confronts his creator. Tyrell asks why he has not come before. “It is not easy to meet one’s Maker,” says Roy. A moment later comes one of my favorite pieces of dialogue. Roy asks whether the maker can change his own creation. “What seems to be the problem?” says Tyrell, like an urbane physician. “Death,” answers Roy. Can Tyrell help? The answer, it turns out, is no, but Tyrell tells him to revel in his time. “The candle that burns twice as bright burns half as long, and you have burned so very, very bright, Roy.” As a consolation, it is dramatically unsuccessful. Though Roy seems to be about to confess to his father—“I have done questionable things”—he then adds, “but nothing the god of biomechanics wouldn’t let you in heaven for.” And on that note, Roy passionately kisses, and then kills, the man who made him. Tyrell, too, has done “questionable things” unrestrained by any god of

biomechanics. The creation is judging the creator. How will we be judged by ours?

Readers find many, many messages in *Do Androids Dream* and in *Blade Runner*—one reason for the enduring power of those works. Clearly, they are both meditations on empathy and personhood, but to me it seems they go beyond a simple injunction that we should love (electronic and genetically altered) others as ourselves. Instead, I see two themes, both central to this book.

First, our beliefs about naturalness, about identity, about empathy, are built on assumptions that start to seem arbitrary, even ludicrous, when presented in the fun-house mirror of an alternative world that is recognizable yet different from our own. We could mock the Voight-Kampff Test for its ironic focus on empathy for animals to deny empathy to androids. What hypocrites and fools the inhabitants of that world are! How blind they are to their own contradictions! But that is not the message I get from *Do Androids Dream* and *Blade Runner*. It is more like “judge not, lest we be judged.” Do you think we will do better? That question should prompt worried humility rather than hubristic condemnation.

The process of empathy that Adam Smith describes may be a basis for ethics, but it will be based on a leap of faith, a projection of identity that will probably rest on a pattern of beliefs full of blind spots and inconsistencies, one that will be in productive tension with our moral theories and our attempts to reason our way to the right answer. Each side will need the other. Our discussions of synthetic personhood will exist in a dialectic between sympathy and syllogism, leaps of empathy and flashes of disgust, hopes that we can realize the angels of our better nature and deep fears that our creations will destroy us. Or judge us. And, as I pointed out earlier, our empathy will be subject to radically conflicting, and partially truthful, primings that cause us to humanize or dehumanize our creations—to push them away into the realm of nonhuman animal, clockwork mannequin, clever software emulation, toaster that can’t say no or, indeed, to recognize them as persons though we know that we made them. What’s more, some of those primings will be designed

into our synthetic creations *on purpose*. Think of the Chimpy deliberately designed to look ape-like and to lack spoken language. Humanities graduate students quickly learn that if they do not know the answer to a question, “it is socially constructed” will generally get them off the hook. Our process of empathy for our synthetic creations will exist in a landscape that is literally made by our own hands. This is “social construction” with a vengeance.

Second, uncertainty about personhood does not stop when we get to our own species line, or our own identity. Deckard is not the only one to doubt his humanity. When we come to explain why *we* are conscious, human, how it is that there is a thinking “I” inside of humanity that has been given dominion over the animal and mechanical world, we could be subject to the same skepticism that the Google engineers directed at Hal. The question is not whether Deckard is a replicant. The question is whether we all are.

Earlier I quoted Turing pointing out that it was hard to prove that humans, too, were more than rule-following automata. He was not the first person to make this argument. In 1887, Samuel Butler had made the same point: “[T]he theory that living beings are conscious machines, can be fought as much and just as little as the theory that machines are unconscious living beings; everything that goes to prove either of these propositions goes just as well to prove the other also.”³⁰ It is (electric) turtles all the way down. To solve the problem of recognizing the Other, it seems, *Do Androids Dream* and *Blade Runner* are telling us we must first know ourselves.

In the movie’s final scene, Roy returns to Sebastian’s apartment alone; we are left to wonder what he has done with the gentle man who helped him and Pris and who tried to flee when he saw Tyrell killed. Roy finds Pris dead, calls her name, cries, daubs himself with her blood, and, howling like a wolf, sets off in pursuit of Deckard. The stroboscope of primplings begins again, flashing like the flickering lights in which the scene is filmed. Roy goes from chanting murderous nursery rhymes, to animalistic howling, to inhuman feats of strength, to Socratic humor, toying with Deckard like a cat with a mouse. “Not very sporting to fire on an unarmed opponent. I thought you were supposed to be good. Aren’t you the . . . *good man?*” Is he?

As the chase continues, we can see that Roy's death is coming closer. His hand seizes up and—metaphor alert—he jams a nail through his palm to keep it open. Yet if Roy is supposed to be a Christ figure, the gift of the god of biomechanics is a horrifyingly ambiguous one, and it is not clear our sins will be expiated by his death.

Finally, Roy has a terrified Deckard defenseless, dangling by one hand from the roof's edge and about to fall. Deckard's fingers slip. Inexplicably, Roy *reaches out*, across that short but enormous divide, grasps Deckard's wrist, and saves him. He gives mercy to the man who has killed his lover when he himself has only moments of life remaining. In those moments, Roy returns to the same point he had raised with Hannibal Chew, the things his eyes have viewed. "I've seen things you people wouldn't believe. Attack ships on fire off the shoulder of Orion. I watched C-beams glitter in the dark near the Tannhäuser Gate. All those moments will be lost in time, like tears in rain. Time to die."

These are profoundly *human* qualities: mercy; the bittersweet confrontation with mortality; the transitory nature of our consciousness and our memory; the impossibility of truly knowing another, of feeling what he or she feels or has lived. The essentially solitary way that all of us, not just the replicants, confront our impending crossing of the line between life and death, person and thing. "At the narrow passage, there is no brother, no friend."³¹ And yet, there is a consolation: our ability through language and art and empathy and dark, dark humor to share something, to reach out across that great divide.

Do we want to accept the gift of Roy's mercy—sparingly dispensed to be sure—but no less remarkable for it? Are we in fact willing to believe what his eyes have seen, to share his memories? Or do we leave him on the other side of the line? All of that, it seems, depends on how we define us as much as how we define him.

Perhaps that is the most important thing to realize from this chapter. Grappling with the question of synthetic Others may bring about a reexamination of the nature of human identity and consciousness that is unparalleled since secular philosophers declared that we would have to learn to live with a God-shaped hole at the center of our world. To draw the line for our creations, we must first draw it for ourselves. We have our own Voight-Kampff Test to face.



2

ARTIFICIAL INTELLIGENCE

There is no security . . . against the ultimate development of mechanical consciousness, in the fact of machines possessing little consciousness now. A mollusc has not much consciousness. Reflect upon the extraordinary advance which machines have made during the last few hundred years, and note how slowly the animal and vegetable kingdoms are advancing. The more highly organised machines are creatures not so much of yesterday, as of the last five minutes, so to speak, in comparison with past time. Assume for the sake of argument that conscious beings have existed for some twenty million years: see what strides machines have made in the last thousand! May not the world last twenty million years longer? If so, what will they not in the end become? Is it not safer to nip the mischief in the bud and to forbid them further progress?¹

THE BUTLERIAN CHALLENGE

That passage was written in 1872. Samuel Butler, the anti-Victorian iconoclast whose novel *The Way of All Flesh* is one of the most searing critiques of the hypocrisies of his time, wrote a book 150 years ago that muses extensively on the possibility of machine consciousness.

Erewhon is a hard book to explain. The title is (nearly) “nowhere” backward—the same thing that “utopia” means in Greek. Erewhon is an imaginary country and it is no utopia. Instead, it is a fun-house mirror in which attentive readers could see Victorian society, and perhaps our society, reflected, reversed. The Erewhonians treat crime the way we do sickness and sickness the way we do crime, imprisoning people for being ill and relying on polite hypocrisies about criminality to excuse their own behavior. How nice it would be to say, “I’d love to come to your party, but I feel some shoplifting coming on.” They punish people for having bad fortune. Arguably, so do we, and that is Butler’s point. Their musical banks

parallel Victorian churches. The currency the musical banks traffic in is honored piously as the true wealth but hypocritically ignored in practice, where real money is what counts. Their universities are “colleges of unreason,” teaching abstruse and archaic doctrines but failing to inspire true critical thinking. As an academic myself, I’ll leave that one alone. Their society even bans the killing of animals and the eating of meat, leading repressed carnivores to feel shame and often contract disease when they finally turn to the black market to gratify their illicit desires. It is a nice parallel to Victorian society’s sexual repression, coupled with its enormous, brutal sex trade. To put it mildly, little in the book is as it seems.

Unwary readers who encounter the two chapters about machine consciousness out of context can be excused for taking them at face value. Was Butler seriously exploring the possibility of machine consciousness? Was he so worried about rogue AI that he even proposed a ban on mechanical progress? Certainly, some people have read him that way. If you know Frank Herbert’s classic science fiction novel *Dune*, you have read about the “Butlerian Jihad” that banned machine intelligences in a distant future. The original Butler would have been amused by that nickname, I think. But just as the musical banks, the courts of illness, and the colleges of unreason are not what they appear to be, the discussion of machine intelligence was mainly supposed to be an allegory for another issue: his era’s passionate debate over the scientific truth and theological implications of biological evolution.

Just what Butler was trying to say is a matter of some dispute. He himself seems either to have been deliberately ambiguous about it or to have changed his position.² Some say he was criticizing evolution, claiming that the same arguments put forward for the gradually increasing complexity of biological beings driven by natural selection would imply that machines could develop consciousness in similar ways. If so, the *reductio ad absurdum* is no longer so *absurdum*. Others say he was using the same form of argument to parody evolution’s critics, and their relentless attempts to suppress, deny, stigmatize, and, if necessary, forbid evolution’s teachings. That one has an unpleasantly modern ring, too.

Butler could have been using machine consciousness as a critical allegory of evolution or an allegory against evolution’s critics. Either way,

a Victorian-era satirical dystopia accurately predicts our contemporary debates about thinking machines. It is as if *Gulliver's Travels* turned out to be a Yelp review of Lilliput as a tourist destination. ("Watch out for the little guys with the ropes! Would rate this place zero stars if I could!") There is a lesson in that. Whether or not he was serious, Butler was right that the same arguments that support biological evolution at least suggest the possibility of machine consciousness. Indeed, as we will see, one possible method of machine learning relies explicitly on an evolutionary mechanism, though the "selfish genes" are algorithms and neural networks running on computers, competing for successful reproduction into the next generation. My imaginary Hal used just such a technique. But Butler is also right that the denunciations of evolution, the explanations of why it is scientifically impossible, will parallel relatively precisely some of the denunciations of AI consciousness and the philosophical explanations that it is impossible. It is worth remembering that the critics were wrong about evolution.

More generally, Butler's work is a good starting place for our discussion for three reasons. First, Butler sees the fragility of the line, its contingent quality. Over the last 40 years, scientists such as the primatologist Frans de Waal have posed skeptical challenges to the idea of a firm, qualitative distinction between humans and nonhuman animals, finding examples of tool use, language, and so on in the animal world. But more than a hundred years earlier, Butler was pointing out that the lines between human and animal *and* human and machine are fuzzier than we might like to imagine. In fact, in words that seem deliberately provocative, Butler challenges both the machine-animal distinction and the idea of qualitatively distinct human consciousness:

Where does consciousness begin, and where end? Who can draw the line? Who can draw any line? Is not everything interwoven with everything? Is not machinery linked with animal life in an infinite variety of ways? The shell of a hen's egg is made of a delicate white ware and is a machine as much as an egg-cup is: the shell is a device for holding the egg, as much as the egg-cup for holding the shell: both are phases of the same function; the hen makes the shell in her inside, but it is pure pottery. She makes her nest outside of herself for convenience' [sic] sake, but the nest is not more of a machine than the egg-shell is. A "machine" is only a "device."³

Having taken a shot at the firmness of the machine-animal distinction, Butler turns to self-awareness. Probably tongue-in-cheek, but no less enlightening for all that, Butler then muses on the consciousness of the humble potato:

Even a potato in a dark cellar has a certain low cunning about him which serves him in excellent stead. He knows perfectly well what he wants and how to get it. He sees the light coming from the cellar window and sends his shoots crawling straight thereto: they will crawl along the floor and up the wall and out at the cellar window; . . . we can imagine him saying, "I will have a tuber here and a tuber there, and I will suck whatsoever advantage I can from all my surroundings. This neighbour I will overshadow, and that I will undermine; and what I can do shall be the limit of what I will do. He that is stronger and better placed than I, shall overcome me and him that is weaker I will overcome." The potato says these things by doing them, which is the best of languages. What is consciousness if this is not consciousness? . . . We find it difficult to sympathise with the emotions of a potato; so we do with those of an oyster. . . . Since . . . they do not annoy us by any expression of pain we call them emotionless; and so *qua* mankind they are; but mankind is not everybody.⁴

Now Butler has the attention not just of the *Dune* reader but the vegetarian, who suddenly realizes that even vegetables might not be fair game. Butler's tongue-in-cheek ode to the possibilities of mind in everything from a steam engine to a potato actually fits into a once-maligned theory of consciousness now enjoying a modest revival. Panpsychism, which dates back to ancient Greece,⁵ claims that mentality or mind is everywhere. It pervades material objects as well as living beings. Adherents run the gamut from mystics to scientists who believe we overstate the differences between animate and inanimate. To be fair, most contemporary panpsychists believe that consciousness reaches its fully developed form only in beings of sufficient complexity, but the potential is there in the humblest of things.

The second reason why Butler is a good starting point for any discussion of the possibility of machine consciousness is even more basic. More than a century ago, he saw that any account of human consciousness that admits it comes from physical interactions in the brain and the nervous system will find it hard to explain why other sets of physical interactions, based on nonorganic processes, cannot produce consciousness. To put it another way, if we deny consciousness to machines because

no true consciousness can come from such a programmed, materialist origin, can we call *ourselves* conscious? Here, again, is Butler from 1887: “[T]he theory that living beings are conscious machines, can be fought as much and just as little as the theory that machines are unconscious living beings; everything that goes to prove either of these propositions goes just as well to prove the other also.”⁶ Seventy years later, Turing would use a similar argument in favor of the Imitation Game, or Turing Test for machine intelligence. If we cannot tell whether an entity is machine or human, even after extensive interaction, who are we to deny another entity consciousness? What ground do we have to stand on?

Finally, Butler’s writing gives me, at least, a timescale for the debate. “The Book of the Machines” was written 150 years ago. The most complex machines around Butler were steam engines, industrial looms, and mechanical calculators. Perhaps one could add the partially completed Babbage Difference Engine, beloved of steampunk science fiction readers and computer science historians. Yet in that context, unimaginably primitive in our terms, he could still see that in the grand sweep of time, “[t]he more highly organised machines are creatures not so much of yesterday, as of the last five minutes.” In other words, he could warn us—with our Siris and ChatGPTs and our deep learning, convolutional neural nets massaging big data—that the timescale of these advances is so short historically, and the pace so rapid, that we should doubt our ability to extrapolate confidently in either direction about the journey’s final destination. That fact should discourage hubris both in those who are skeptical Artificial Intelligence will ever be developed, and those who are confident that it will arrive in some specific anticipated format and revolutionize the world in the very near future. Hubris, however, appears to be an endlessly renewable resource.

HUBRIS AND HUMILITY IN AI

The history of AI is a history of overconfident predictions. In August 1955 a group of academic luminaries submitted a grant proposal to the Rockefeller Foundation for a summer workshop on AI. The document is famous partly for its historical importance—and it is a grant proposal. Every time I read it, I find myself imagining equivalent texts from other

historical moments. (“Executive Summary: Goal: to escape from slavery under Pharaoh. Needs: Method of parting the Red Sea. Also, snacks.”) But the document is also famous for its ambition—beginning a dialectic in AI research between wildly optimistic claims and pessimistic laments of difficulty that continues to this day. Note the goals:

We propose that a 2 month, 10 man study of artificial intelligence be carried out during the summer of 1956 at Dartmouth College in Hanover, New Hampshire. The study is to proceed on the basis of the conjecture that every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it. An attempt will be made to find how to make machines use language, form abstractions and concepts, solve kinds of problems now reserved for humans, and improve themselves. We think that a significant advance can be made in one or more of these problems if a carefully selected group of scientists work on it together for a summer.⁷

For a *summer*. Progress was not quite as fast as they imagined. Nevertheless, ten years later, giants in the field such as Marvin Minsky and Herbert Simon were predicting General-Purpose Artificial Intelligence or “machines . . . capable . . . of doing any work a man can do” by the 1980s.⁸ Huge strides have been made in aspects of artificial intelligence—machine-aided translation, facial recognition, autonomous locomotion, expert systems, and so on. But General AI—an intelligence that exhibits all the qualities of human intelligence and capability—has remained out of reach. Indeed, because the payoff from these more limited subsystems—which today power everything from Google Translate and image recognition to the recommendations of your streaming service—is so rich, some researchers have argued that the goal of General AI was a snare and a delusion. What was needed instead, they claimed, was a set of ever more powerful subspecialties—expert systems capable of performing discrete tasks extremely well but without the larger goal of achieving consciousness or passing the Turing Test. There might be “machines capable of doing any work a man can do,” but they would be multiple different machines, with no ghost in the gears, no claim to a holistic consciousness.

It is worth noting that, under some definitions, that might be enough to be hailed as Artificial General Intelligence. For example, Metaculus, a site that solicits and aggregates predictions of future events, has as its criteria for high-level General AI that it has to be able to pass a two-hour adversarial Turing Test featuring text and images, assemble a complex

model car, and perform well on tests assessing a number of other capabilities. The focus is on capabilities. In other words, if we could have a machine that did all of the things humans can do, from composing a sonnet to conversing fluently, from changing a lightbulb to piloting a plane, that would be enough. The development of such a multitalented machine would certainly transform our economy and society, but my interest is in AI *personhood* and potential consciousness; being an extremely competent collection of expert systems is not automatically enough. Beyond those skills, I am asking the question of whether there is some consciousness, some set of morally salient capabilities, that would cause us to see the machine as a moral actor whose personhood should be recognized.

Despite the history of overconfidence and of setbacks, arguments that General AI will appear in the near future have not ended. Indeed, if anything, the optimistic claims have become even more far-reaching. Thirty years ago the buzzword among the most fervent AI optimists was the Singularity, a sort of technological liftoff point in which a combination of scientific and technical breakthroughs lead to an explosion of self-improving Artificial Intelligence coupled to a vastly improved ability to manipulate both our bodies and the external world through nanotechnology and genetic engineering.⁹ Writers such as Vernor Vinge and Ray Kurzweil used the term Singularity to refer to the point where, because of exponential technological growth, the graph of technological progress will go vertical or at least be impossible to predict using current tools. Assuming explosive and imminent advances in AI, they believed that we would soon have improvements not in technology alone, but in the intelligence that will create new technology. Intelligence itself will be transformed. Once we have built machines smarter than ourselves—machines capable of building machines smarter than themselves—we will, by definition, be unable to predict the line that progress will take. Vinge, whose 1993 article¹⁰ initiated the focus on an AI Singularity, was pessimistic about what might result. Why should we assume that an intelligence vastly greater than our own would treat us any better than we treat chimpanzees? Kurzweil, by contrast, generally saw the Singularity leading us into a glorious world of posthuman immortality.

Kurzweil's view seemed to resonate more in frothy, popular science discussions, but, in recent years, an alternative to Kurzweil's view has

developed, one that hearkens back to Vinge's original caution. This perspective, associated with researchers such as Eliezer Yudkowsky and Nick Bostrom, shares with Kurzweil the intuition that Artificial General Intelligence may arrive much sooner than many of us expect. It differs in that the consequences it foresees are by no means as benign. In Yudkowsky's words, "The AI does not hate you, nor does it love you, but you are made out of atoms which it can use for something else."¹¹ Far from bringing us immortality and a peaceful and insanely productive, ecologically sustainable world, Yudkowsky and his fellow skeptics argue that the Singularity could bring global devastation and even human extinction.

The term "singularity" is actually drawn from a memorial tribute given by Stanisław Ulam to the famous mathematician and information theorist John von Neumann.¹² It is normally quoted in an abbreviated form that suggests von Neumann's eminence can be enlisted in support of the optimistic Singularity vision. Read in full and in context, however, the original quotation uses the term "singularity" to refer to a different and less positive set of possibilities than Kurzweil's image. Ulam says of von Neumann:

Quite aware that the criteria of value in mathematical work are, to some extent, purely aesthetic, he once expressed an apprehension that the values put on abstract scientific achievement in our present civilization might diminish: "The interests of humanity may change, the present curiosities in science may cease, and entirely different things may occupy the human mind in the future." One conversation centered on the ever accelerating progress of technology and changes in the mode of human life, which gives the appearance of approaching some essential singularity in the history of the race beyond which human affairs, as we know them, could not continue.¹³

Far from racing with delirious optimism into a technologically transformed future, I read von Neumann, and perhaps Ulam, to have apprehensions about the "changes in the mode of human life" in a future that they could not predict and in which "human affairs, as we know them, could not continue." This is hardly the full-throated endorsement of the optimistic Singularity. In fact, it sounds a Burkean note of caution that would later be echoed in Bostrom and Yudkowsky's darker visions of how AI might transform or destroy our world.

On the other hand, von Neumann is putting forward two premises central to the contemporary usage of the term. First, technological

progress—or at least technological progress in some fields—is exponential, not linear. (But for how long?) Second, while the first few stages of an exponential graph are not that different from a linear one, the line on the graph quickly goes almost vertical. This will lead those who are assuming more linear growth, or who are standing on the flatter part of the time curve, to dramatically overestimate how long technological developments will take to achieve. It will also rapidly put the future out of sight from where we are, thus rendering it impossible to predict. Strikingly, despite this fact, some of the proponents of the Singularity *do* prophesy with apparent confidence about what will transpire after it occurs. Kurzweil imagines a posthuman, technologically enabled immortality, for example.

To the uninitiated, the future painted in Kurzweil's 2005 *The Singularity Is Near* sounds like a delightfully wacky fantasy, a high-tech version of the rapture in which our posthuman bodies rise up to an endless virtual reality in the cloud, run by benign intelligences that have long ago transcended our limits. A “version” of the rapture? That *is* the rapture. No wonder the more enthusiastic odes to the Singularity have a religious, chiliastic feel to them. Sometimes, that impression can get in the way of a careful assessment of the specific claims being made about AI that, while overly optimistic, are based on thought-provoking premises.

If technological change (e.g., the doubling of computer chip capacity every 18 months to two years that is known as Moore's law) could continue on an exponential curve, then a dramatically different future will arrive far sooner than we expect. That is Kurzweil's central point, as it was von Neumann's. But many scientists warn that we are rapidly approaching the physical limits of science in making transistors smaller. What's more, some have argued that, at our current levels of technology, cost-benefit analysis will no longer support the titanic investments required to continue to meet that benchmark. Moore's law may have ceased to be true already. The exponential graph may flatten out, whether it is flattened by physics or balance sheets or both.

To be fair to those who believe in a short timeline to General AI, they generally do not predict a single, invariant, exponential curve but rather a stacked series of S-curves in which a particular technology starts off slowly, hits an exponential period of innovation, flattens off, and is in

turn replaced by a new technology that goes through the same stages. One way for this trend to continue in the realm of computer architecture would require us to predict, for example, that current chip designs would be overtaken by a new paradigm—quantum computing, say, which would exploit the physics of the quanta such as the entanglement of quantum particles, Einstein’s “spooky action at a distance.” Of course, the dramatic advance does not have to be quantum computing. Perhaps Richard Feynman was right and there is still room at the bottom, in the nanoscale, using technologies and heat dissipation methods we only dimly understand now. Or perhaps some combination of biological computing and machine computing will open the next frontier. Perhaps the transformation will not primarily be to the hardware at all but rather in the software, with new techniques of machine learning producing quantum leaps in performance. Regardless of the specific technique, the large claim is that we will continue to find new revolutionary technologies that will enable yet another S-curve in computer capacity. Yet how can we confidently predict such paradigm shifts in technology? By definition, they are outside of our current technological frame of reference.

The speed of technological transformation will be particularly hard to predict if we are talking about *multiple* technologies, sometimes accelerating on exponential curves, having unexpected synchronistic effects on each other. Take the evolution of computer networks from 1990 to 2005, for example. Most of the basic technological components of the internet were there in the 1980s. Versions of the internet itself—a distributed packet-switched system—date back to the 1950s. But during this period of time, those things suddenly came together to form the World Wide Web, to revolutionize our communications, our media, and our global commerce.

We can debate what addition supersaturated the solution and precipitated the crystal of transformative innovation—Tim Berners-Lee’s architecture of HTML and the World Wide Web? The price, speed, and memory frontier that PCs hit in the early 1990s? The unused bandwidth available on cable networks’ fiber backbones due to networks and “rights of way” property regimes created for an entirely different purpose? More likely, it is all of the above. Without any single great breakthrough, the world was suddenly dramatically different. The worldwide internet went from

being a science fiction trope that was never going to exist (“flying cars!”) to a reality in about five years. It became *the* reality—an unquestioned feature of our world like gravity and oxygen—in a mere 15. People who had predicted for decades that computers and networks would transform society, and faced entirely justified heckling when the promised revolution failed to appear, were wrong, wrong, wrong until they were suddenly and shockingly right. In 15 years, the world changed dramatically, without warning and without some eureka discovery that might have been thought necessary to precipitate the transformation. All the technologies were well understood. The result was not. That incident is undeniably part of our past. And we think we can predict the future?

Why can this not happen with AI? I do not mean to say that it will, but confident assertions either way should be met with skepticism. Duke Law School’s parking lot has some gratifyingly witty bumper stickers. One seems appropriate here. “Radical Agnostic” says the large, capitalized text. Underneath is the smaller punch line. “I don’t know and *you don’t either!!*” Perhaps this should be our motto for AI prognostication. Some may think, perhaps rightly, that I fail that test. I am going to argue that there are reasons to believe that progress is likely to be faster than many of us think. My agnosticism has a tilt. Nevertheless, I think the radical agnostic’s motto is the right one.

If the internet’s transformation seems too singular and unlikely to be representative, it is worth remembering that we have just lived through another example of this process of synchronistic change: the rapid proliferation of neural network systems that rely on deep learning to recognize speech in multiple languages, translate sentences, identify pictures, predict consumer desires, and so on. How did this happen? The origins of electronic neural networks can be found as far back as the 1940s and 1950s.¹⁴ A cluster of events had to come together to produce the leap forward of the last ten years. There were revolutionary breakthroughs in network theory and design—the software side. Continuous improvements in speed and drops in cost of hardware made those software advances suddenly have a much greater reach or potential. But wider cultural and technological transformations also played a role. Both the software and hardware showed what they could do because of an explosion of data on which they could be tested and proven.

Combine the continuously improving technologies of the individual computer—which is rapidly increasing in speed, processing power, and memory capacity while dropping in price—with a global network of other computers doing the same thing and a cloud that is almost always in reach. Put those computers in people’s pockets, as smartphones. Now we have nearly seven billion nodes connecting to the cloud around the planet, each performing a host of different tasks and running many different apps, and thus an *exponential increase* in the rate of data generation by those rapidly proliferating devices.

Millions of people navigate using Google Maps, upload and tag photographs, dictate commands to their phones, and then correct that dictation, providing feedback to the system. The torrent of data is staggering—“big data,” indeed. And in that data are patterns, patterns that artificial intelligence can “learn” to identify. Rather than programming the system with rules up-front—“this is the shape of a cat,” “when a British person says ‘bath’ it sounds like this”—the system uses an architecture very loosely emulating the organization of neurons in the brain, arranged in sequential layers of processing. Many such systems develop through a process of trial and error, giving greater weight to the input from those layers that improve the accuracy of predictions. Once programmed with goals and parameters, and in some cases with an initial curated data set, the system can perform this process on its own, layer after layer, developing its own credit-assignment paths that lead to ever more precise identification in a process that may be partially inscrutable even to the original programmers. The system might even be given almost no guidance and simply rewarded through deep reinforcement learning when it does something its programmers think is good. This technique has consistently outperformed more structured, choreographed approaches to the problems machine intelligence must solve.

Look at the number of technological developments that come together to make this happen. It is not simply a matter of Moore’s law, which skeptics rightly point out is no longer empirically accurate. Deep learning depends on dramatic changes in memory capacity, price, distributed storage, number of users, and advances in artificial intelligence theory and software. It turns those advances onto the firehose of data generated by our computer systems. And the neural network uses deep learning,

rather than some formal set of preprogrammed rules, to master this torrent of data. Peter Norvig, the Director of Research at Google and a leading scholar of artificial intelligence, puts it nicely: “We decided that the best model of the world was the world.”¹⁵

Deep learning has been a revolutionary development.¹⁶ Google Translate became dramatically better literally overnight. Image- or speech-recognition software was suddenly vastly more accurate. For all of this, you have deep learning, and probably neural networks, to thank. What does this tell us about the prospects of General AI? By itself, not much at all. True, this is one type of artificial intelligence, focused on discrete tasks, but it is not General AI, let alone consciousness, unless your threshold for consciousness is “can you identify all the cute little kitty cats in this picture?” Large language models such as ChatGPT or LaMDA are such systems. Blake Lemoine, the Google engineer whose story began this book, was so convinced by LaMDA’s output that he believed it had become conscious. Lemoine was incorrect: there is no ghost in that machine, merely jaw-droppingly brilliant imposture.

The story of deep learning, and of Lemoine’s error, do not teach us that General AI is here, or that machine learning systems like LaMDA or ChatGPT are going to become conscious tomorrow. Instead, they should teach us something very different: that it is very hard to forecast developments in technologies, some of which are developing at exponential rates, when it is the *interaction* of the rapidly changing components of the system that enables the dramatic, paradigm-shifting change. The point is that sudden and unexpected change is possible, though not inevitable, whether from exponential growth within one field or syncretic fusion among many. That suggests we might want to take seriously the arguments of those who think that Artificial General Intelligence may arrive sooner than we think, even if we are skeptical of their precise timetable, or their predictions of rapturous immortality or machine-led annihilation. We need not rely on their arguments as descriptions of what *will* happen and *when*. We can think of them as reasonable suggestions of what *could* happen, and *why*.

Perhaps an anecdote will underline that point. As I was writing these words, I saw the news that Geoffrey Hinton, a renowned pioneer in neural networks, had resigned from Google so that he could speak more freely

about his concerns over AI systems. This was not exactly like Thomas Edison quitting his job because he was worried about the effects of light-bulbs but, for many in the field, it produced an equivalent level of shock. To be clear, Hinton's concerns about the breakneck pace of technological development around AI systems are broad ones. He instanced everything from the rampant production of deep fakes to the effect on the labor market, warfare, and political stability. But I was struck by one thing he said: "The idea that this stuff could actually get smarter than people—a few people believed that. . . . But most people thought it was way off. And I thought it was way off. I thought it was 30 to 50 years or even longer away. Obviously, I no longer think that."¹⁷ He is not alone in this belief. Google's DeepMind is without doubt one of the most important companies in the field. Its research has been vital to current breakthroughs, including those by rival companies. The same month as Hinton's resignation, DeepMind's CEO Demis Hassabis had this to say: "The progress in the last few years has been pretty incredible, I don't see any reason why that progress is going to slow down. I think it may even accelerate. So, I think we could be just a few years, maybe within a decade, away [from human-level AI]."¹⁸

ARTIFICIAL INTELLIGENCE? WHEN?

This brings us to the obvious question: Will General-Purpose, or even conscious, AI arrive at all, and if so, when? It turns out that those studying AI have radically different answers to those questions. They differ about the most promising lines of research, their difficulty, and the extent to which industry and academic research scientists will actually be focused on Artificial General Intelligence rather than on building many discrete artificial intelligence systems that make hair appointments, book your travel, or organize your photo album. But they also differ on the two axes just identified: optimism or pessimism about sustained exponential growth and optimism or pessimism about the frequency and significance of technological synchronicity—the coming together of many factors to produce a leap forward that was not predictable in advance.

These forms of optimism and pessimism are shared in the discussion of economic growth more generally, of course. Tyler Cowen's *The Great*

*Stagnation*¹⁹ and Robert Gordon's *The Rise and Fall of American Growth*²⁰ both provide compelling arguments against the assumption that we will continue to have the kind of robust economic growth, year after year, that characterized much of the twentieth century, though Cowen is actually more optimistic.²¹ But the AI debates present a particularly hard puzzle for prediction because we have glaring examples of remarkable, and in some cases exponential, rates of technological advance. Yet we also have repeated, humility-inducing difficulties and failures. After all, some problems that AI scientists at first thought were fairly basic (teaching a computer "common sense," for example) have proven remarkably hard to solve:

A.I. "recognizes objects, but can't explain what it sees. It can't read a textbook and understand the questions in the back of the book," said Oren Etzioni, a former University of Washington professor who oversees the Allen Institute for Artificial Intelligence. "It is devoid of common sense." Success may require years or even decades of work—if it comes at all. Others have tried to digitize common sense, and the task has always proved too large. In the mid-1980s, Doug Lenat, a former Stanford University professor, with backing from the government and several of the country's largest tech companies, started a project called Cyc. He and his team of researchers worked to codify all the simple truths that we learn as children, from "you can't be in two places at the same time" to "when drinking from a cup, hold the open end up." Thirty years later, Mr. Lenat and his team are still at work on this "common sense engine"—with no end in sight.²²

That skepticism could be strengthened by a series of disagreements in the field about the best methods for developing even discrete expert systems, let alone Artificial General Intelligence. Should AI be neat or scruffy? Neat approaches are based on some overarching framework such as symbolic logic, and they use that framework to solve every problem. Scruffy approaches, by contrast, opportunistically use different cognitive techniques to solve different problems so that the method for translating from one language to another might be different than the method for image recognition or playing chess, and much might consist of ad hoc, individually coded heuristics based on real-world experience. Should or will AI be rule governed, based on an enormously complex but finite set of algorithms laid down at the start by its designers? Alternatively, will it be partially autonomous, "learning" how to achieve tasks in ways that may be inscrutable to the original creators? Will it be based on advances

in the logical dissection of how humans actually think or on the pursuit of rational problem-solving, regardless of how humans think? Something else altogether? If the AI optimists cannot even tell us what methods will yield General AI, then how can their optimism be sustained?

Reflecting the number of questions to be answered, surveys of AI researchers have shown considerable divergence in predictions of when General AI, or something like it, would be achieved. One notable 2016 survey²³ used as its target population all of the researchers who published at two of the most important conferences in the field and asked, among other things, when high-level machine intelligence would be achieved. Their definition of such intelligence was a demanding one: “High-level machine intelligence (HLMI) is achieved when unaided machines can accomplish every task better and more cheaply than human workers.”²⁴ Note that this definition, like any we might choose, will have dramatic effects on the outcomes. For example, we might want to know when the first example of General AI could be achieved if we were willing to put Manhattan Project-level resources into it, not when every doctor, novelist, lawyer, composer, and kindergarten teacher could be replaced by a better, and cheaper, cybernetic equivalent. Alternatively, if our concerns were with the question of when there might be some *moral* claim to legal personhood, we might think it irrelevant whether the AI could do brain surgery or dance ballet, just as long as we felt its consciousness shared enough with our own to warrant such a claim. The advantage of the question the researchers posed is that it looks formalizable and falsifiable, avoiding philosophical debates about whether true consciousness had been or ever could be achieved. That is also its disadvantage. Still, given both its universality of field and its price constraint—*every* task humans can do, in *every* case, done cheaper—it presents a very demanding standard.

The aggregate forecast was that there was about a 30 percent chance of achieving high-level machine intelligence within about 25 years (as of 2016) and a 50 percent chance of achieving it within 45 years. The researchers reported a striking demographic split in responses: “Asian respondents expect HLMI within 30 years, whereas North Americans expect it in 74 years.” Interestingly, the aggregate forecast suggested there was a 10 percent chance that it might be achieved within nine years of

2016, that is by 2025! Kurzweil's view is still an outlier, but it falls, or fell, within the bounds of the profession.

As of August 2022, Metaculus, the online prediction site, was forecasting that we would have Artificial General Intelligence by November 2041.²⁵ Their criteria for Artificial General Intelligence were different than the survey above; as I mentioned before, the system had to be able to perform well on tests assessing varied skills ranging from a two-hour adversarial Turing Test, featuring text and images, to the assembly of a complex model car. By May 2023, their assessment had changed. "The Metaculus community currently expects [Artificial General Intelligence] to be unveiled in October 2031." The influential AI thinker Eliezer Yudkowsky showed equal optimism about the speed of the transformation, coupled with extreme pessimism about its results. He accepted the following bet from Bryan Caplan: "Bryan Caplan pays Eliezer \$100 now, in exchange for \$200 CPI-adjusted from Eliezer if the world has not been ended by nonaligned AI before 12:00am GMT on January 1st, 2030."²⁶

On the other end of the spectrum from the Singularists are skeptics who find these predictions wildly optimistic (or pessimistic, depending on what you think General AI will do when it arrives). Rodney Brooks, a former director of the MIT Computer Science and Artificial Intelligence Laboratory, and the founder of iRobot, the company that makes your Roomba, has been a frequent critic of overconfident predictions. He claims they are characterized by a pattern of fallacies. They predict consistent exponential rates of technological growth rather than a regression to the mean. They use trivial accomplishments (iPhoto recognizing all the photos of your lover's face) as evidence for the idea that qualitative transformations (General AI) are close at hand. Finally, they make firm technological projections when the timescale means that neither the technology nor the state of the world in which that technology will be deployed can accurately be predicted.²⁷ Brooks pointedly rejects Kurzweil's claims, and some of his own projections put human-level AI much further in the future: "It will be well over 100 years before we see this level in our machines. Maybe many hundred years."²⁸ Interestingly, though, it is the optimistic time-horizon and suddenness suggested by the proponents of the Singularity that Brooks doubts, not the eventual achievement itself. Instead, he imagines a gradual process of improvement, "generation by

generation by generation. The singularity will be a period, not an event.” We will be driven, he thinks, “not by the imperative of the singularity itself but by the usual economic and sociological forces. Eventually, we will create truly artificial intelligences, with cognition and consciousness recognizably similar to our own.”²⁹

Why is Brooks so confident, given that he is generally a skeptic of optimistic AI claims? The reason is simple. We are learning more and more about the neurological processes of the brain. What we can understand, we can hope eventually to replicate:

I, you, our family, friends, and dogs—we all are machines. We are really sophisticated machines made up of billions and billions of biomolecules that interact according to well-defined, though not completely known, rules deriving from physics and chemistry. The biomolecular interactions taking place inside our heads give rise to our intellect, our feelings, our sense of self. Accepting this hypothesis opens up a remarkable possibility. If we really are machines and if—this is a big if—we learn the rules governing our brains, then in principle there’s no reason why we shouldn’t be able to replicate those rules in, say, silicon and steel. I believe our creation would exhibit genuine human-level intelligence, emotions, and even consciousness.³⁰

This is not the most likely method of achieving General AI, far from it. Think of Brooks’s postulate as an upper bound in AI research—one way of conceiving of the problem that indicates General AI *must be* achievable, if incredibly hard. We have a model of a functioning consciousness: us.

Some will believe that, by divine command, consciousness can only be created by the deity, not by human hands and minds. Perhaps there is some as-yet-undiscovered emergent property of natural biological brains that cannot be reproduced, even if replicated perfectly, either *in silico* or even in some biological computational device. Others believe that consciousness is, in some strange way, *prior* to material reality—the substrate on which the observable physical universe depends—though this still begs the question of whether machines could have the requisite consciousness. But barring a divine or technologically intractable limit—some neurological equivalent of the light-speed barrier—eventually we will be able to recreate the relevant aspects of our brains and hence our consciousness. Having done that, we might be able to transcend some of the human brain’s limitations in terms of speed, memory capacity,

embedded knowledge base, and networked communication of thought. Starting with a model based on a physical brain we could create ever more capable forms of general, conscious Artificial Intelligence. This is extremely unlikely to be the way we would achieve General AI. In fact, it might be the hardest and the one that would take the most time. But reconceived this way, the problem becomes a material and a soluble one. And Brooks, remember, is a skeptic.

IT'S ALL ABOUT THE HARDWARE(?)

Writers on AI agree that neither the range of predictions nor the fact that the due date keeps getting bumped forward induce confidence. In his seminal 1993 article, Vinge acknowledges this fact when making his own prediction. “I believe that the creation of greater than human intelligence will occur during the next thirty years. (Charles Platt has pointed out that AI enthusiasts have been making claims like this for the last thirty years. Just so I’m not guilty of a relative-time ambiguity, let me more specific: I’ll be surprised if this event occurs before 2005 or after 2030.)”³¹ This aside became known as Platt’s Law: those making predictions about General AI will place its inception date roughly 30 years in the future from the date the prediction was made.

Is there some less-subjective basis on which we could predict General AI? Are there metrics that would provide us a benchmark for progress? One answer is that we do not need to replicate the specific architecture of the brain but rather to emulate, in silicon or its successors, all of the relevant capacities and capabilities of a brain—the amount of memory it can hold, how fast it can solve problems, and so on. (Hal, the imaginary computer from the introduction, achieved sentience when the number of connections in his neural networks hit a number that approximated that of a human brain. But that was a thought experiment. There is no reason to think this is the relevant metric.) Once we have equivalent hardware, goes the theory, we only need to tweak the software, and voila, General AI! But where are we in terms of comparative capabilities? And what is the historical rate of change? In 2011, eons ago in internet time, *Scientific American* ran the article “Computers versus Brains”:

For decades computer scientists have strived to build machines that can calculate faster than the human brain and store more information. The contraptions have won. The world's most powerful supercomputer, the K from Fujitsu, computes four times faster and holds 10 times as much data. And of course, many more bits are coursing through the Internet at any moment. Yet the Internet's servers worldwide would fill a small city, and the K sucks up enough electricity to power 10,000 homes. The incredibly efficient brain consumes less juice than a dim lightbulb and fits nicely inside our head. Biology does a lot with a little: the human genome, which grows our body and directs us through years of complex life, requires less data than a laptop operating system. Even a cat's brain smokes the newest iPad—1,000 times more data storage and a million times quicker to act on it.³²

All of these figures, except those claimed for the brain, which are problematic for other reasons, are now out of date, of course. The 2011 *Scientific American* article claims that the K supercomputer could then perform 8.2 petaflops or 8.2 *quadrillion* (8.2×10^{15}) floating-point operations per second. That was a marked advance from earlier computers. As late as 2008, IBM's Blue Gene, the fastest supercomputer at the time, was just above 1 petaflops. By contrast, the Frontier, the fastest supercomputer as of 2023, can perform 1194 petaflops, 145 times faster than the K and 1100 times faster than the Blue Gene. From the Blue Gene to the Frontier, processing speed doubled approximately every 18 months. While this may not exactly be exponential growth, it is a startling rate of improvement. And this comparative hardware approach leads people other than proponents of the Singularity to be fairly optimistic about how soon General AI will arrive. To quote Nick Bostrom, the Oxford University professor whose book *Superintelligence: Paths, Dangers, Strategies* warns of the dangers rather than the promise of AI:

Hardware-wise, the brain still compares favorably with machines. Estimates vary, but perhaps the cortex performs something like 10^{16} or 10^{18} operations per second using 20 watts, which is impressive. Eventually, the limits of computation in machine substrate are of course far beyond those in biological tissue, and it shouldn't take too long to reach rough equivalence. The advance of algorithms is harder to predict, but the notion that we could have human-level AI within a small number of decades seems credible, though there is great uncertainty on both the lower and upper sides of this estimate.³³

Bostrom's estimate of the brain's capacity is higher than that of the *Scientific American* article. The authors of that piece estimated the brain

could perform 2 petaflops. Bostrom seems to assume that it can perform somewhere between 10 and 1,000. Jürgen Schmidhuber, scientific director of a leading Swiss AI Lab and a machine learning pioneer, is also optimistic about the arrival of General AI. His optimism is based not just on the absolute speed of the very fastest machines but on the falling price of the average machine:

When will we have computers as capable as the brain? Soon. Every five years computing is getting roughly 10 times cheaper. Unlike Moore's Law, which says that the number of transistors per microchip doubles every 18 months (and which recently broke) this older trend has held since Konrad Zuse built the first working program-controlled computer. His machine could perform roughly one floating-point operation per second. Today, 75 years later, hardware is roughly a million billion times faster per unit price. Soon we'll have cheap devices with the raw computational power of a human brain; a few decades later, of all 10 billion human brains together, which collectively probably cannot execute more than 10^{30} meaningful elementary operations per second.³⁴

The Open Philanthropy Project, an effective altruism nonprofit, has funded a lot of research on the possible impact of AI. In 2020 they commissioned a report on when we might have human-level AI. The author of that report,³⁵ Ajeya Cotra, found a 10 percent chance by 2031, a 50 percent chance by 2052, and an almost 80 percent chance by 2100. She used a number of methods, including the "floating-point operations in the brain" analysis we have just been discussing. She even attempted, as one benchmark, to estimate the number of floating-point operations represented by the entire history of biological evolution toward humans. It is as if we saw biological evolution as a moonshot AI project trying to achieve human consciousness and could extrapolate from that how long it would take machines affordably to replicate that evolutionary path. Cotra then adjusted the sum of all of these predictive models and the median fell on, specifically, 2052, or 32 years after the report was published. A cynic might say that Platt's Law still holds! Two years later, Cotra adjusted her median prediction to 2040 because of unexpectedly good performance on a number of benchmarks since 2020.³⁶

But what do all these numbers actually mean? A critic might say that they are fundamentally misleading. Human beings do not think in floating-point operations. You can indeed calculate 1.37×8.91 , but I am fairly sure you don't do it in a single second, still less in a millionth

or billionth of a second. Nor do we conceive of the activities of recognizing a face, realizing your marinade needs more ponzu, or writing a love poem as involving floating-point operations at all. Is using this number to compare the power of a brain and a computer like using miles per hour to quantify Shakespeare's prose? To paraphrase Norvig and Russell's book *Artificial Intelligence*, we do not compare the albatross and the 747 by asking how quickly each flaps its wings. They achieve flight using different techniques and, barring the attempt to replicate the brain neuron by neuron, the same will be true of an attempt at building General AI.

Are these comparisons useless, then? Despite the criticisms I just pointed out, as long as they are taken with an appropriate degree of caution, such comparisons do help illuminate something useful. Any attempt to create General AI is aided by having more capable, faster, cheaper, smaller computers, which can handle more complex sets of instructions, contain more memory, form networks more easily, and so on.

In the past, artificial intelligence researchers have found that increases in speed mean problems that were once thought to require elegant solutions may in fact be solved by brute-force approaches. For example, we might think the only way to teach a computer to play chess is by elaborately programming software rules that outline strategy and tactics. Or perhaps just to have the computer teach itself by playing millions or billions of games, generating its own rules and strategies, using a technique called deep reinforcement learning. When I interviewed Hal Abelson, a renowned computer scientist at MIT, he told me that "problems that people thought could only be solved elegantly are instead being solved by simple techniques of reinforcement learning."

One of the most powerful examples of reinforcement learning is provided by the development of DeepMind's Go-playing system. The game of Go has vastly more permutations than chess: "As simple as the rules may seem, Go is profoundly complex. There are an astonishing 10 to the power of 170 possible board configurations—more than the number of atoms in the known universe. This makes the game of Go a googol times more complex than chess."³⁷ With a game this mind-numbingly complicated, it would seem that any AI would have to emulate human strategies of intuition and pattern recognition and would have to rely

on the tactical heuristics polished by generations of players—or not. The researchers at Google’s DeepMind project created a program called AlphaGo, which went on to beat the best human players in the world. The first version of AlphaGo was “trained by supervised learning from human expert moves, and by reinforcement learning from self-play.” These techniques rely on a curated dataset and an initially supervised interaction with that dataset. That is still far less direction, far less programmed strategy, than researchers had previously believed would be necessary. Yet the results of its victorious contests with human grand masters were remarkable: “During the games, AlphaGo played a handful of highly inventive winning moves, several of which—including move 37 in game two—were so surprising they overturned hundreds of years of received wisdom, and have since been examined extensively by players of all levels. In the course of winning, AlphaGo somehow taught the world completely new knowledge about perhaps the most studied and contemplated game in history.”

To find the limits of deep reinforcement learning, the researchers created a second version of the program, called AlphaGo Zero,

based solely on reinforcement learning, without human data, guidance or domain knowledge beyond game rules. AlphaGo becomes its own teacher: a neural network is trained to predict AlphaGo’s own move selections and also the winner of AlphaGo’s games. This neural network improves the strength of the tree search, resulting in higher quality move selection and stronger self-play in the next iteration. Starting *tabula rasa*, our new program AlphaGo Zero achieved superhuman performance, winning 100–0 against the previously published, champion-defeating AlphaGo.³⁸

In the words of the AlphaGo Zero team: “This technique [of reinforcement learning without human guidance] is more powerful than previous versions of AlphaGo because it is no longer constrained by the limits of human knowledge.”³⁹

To be clear, AlphaGo Zero is not General AI or anything remotely close to it. It also was not achieved solely because of hardware advances; the researchers at DeepMind are justifiably proud of their astonishing accomplishment in both software and neural architecture design. ChatGPT would not exist without those advances, which the DeepMind teams

shared widely. But the increase in speed, memory, and data-handling capacity that I described earlier opens entirely new possible lines of research. Neural networks, deep learning, and reinforcement learning show that we can achieve striking results at tasks previously thought to play to human strengths, such as facial recognition or intuitive strategy games, without attempting precisely to emulate the patterns of human thought.

What does this tell us? There is no one-to-one map of human and machine capability; at least at the moment, both the hardware and software are very different. Thus, the head-to-CPU comparisons of processing capabilities are wildly approximate at best. But if one trims away the hyperbole about operations per second, and number of neural connections, a truth remains. While we do not know what the crucial dimensions of hardware performance will be in eventually achieving General AI, the rate of progress on *every* dimension of performance suggests that Bostrom and Schmidhuber have reason for their qualified optimism. As with military strategy, greater resources mean more angles of attack, some of them previously unforeseen.

Kurzweil, of course, believes that General AI is much closer: “When will we have computers as capable as the brain? I believe computers will match and then quickly exceed human capabilities in the areas where humans are still superior today by 2029.”⁴⁰ Yudkowsky, in his pessimism, seems to believe that there is a significant danger of us achieving General AI not long after that date. From my discussions with AI researchers, I find this prediction unlikely, though some of them have become decidedly more optimistic recently. But I find equally puzzling those who claim confidently we are centuries away. The graph of technological change may not be vertical, but it is steep and punctuated by unforeseen leaps forward, sometimes driven by the synchronicity of multiple technologies unexpectedly coming together, sometimes by new approaches that harness rapidly evolving speed and big data capabilities, and sometimes by theoretical breakthroughs. At the very least, I think we can be confident of this: long before the century is out, we will have AI at a level where its consciousness is, at least, a matter on which well-informed people can, and will, reasonably disagree. The controversy will be live. Indeed, some would argue we are already there. And that is all I need.

EVEN IF IT WORKS, IS IT CONSCIOUS?

If one challenge to General AI is that it is impossible, or will take hundreds of years to achieve, a second and more fundamental challenge goes to *ontology* rather than *technology*: the nature of being, not the likelihood of working. Even if a computer-based Artificial Intelligence could do anything a human could do, would we think it was alive, aware, and thus perhaps a person? After all, it is just a machine. It is doing only what it has been programmed to do. It might replicate our responses with perfect fidelity, but would it be conscious while doing so or merely parrotting lines programmed by others, like Siri “remembering” your birthday and congratulating you on it? Let us begin with Alan Turing and his critics.

In “Computing Machinery and Intelligence,”⁴¹ Turing poses the question, “can machines think”? He then quickly suggests substituting for that question, which he calls “meaningless,” another one: Can an interrogator distinguish between a human being and a machine on the basis of their typed answers to the interrogator’s questions? Turing’s reasons for proposing this substitution are not exactly clear. He says that it “has the advantage of drawing a fairly sharp line between the physical and the intellectual capacities of a man.” He says that one alternative method of answering the question “can machines think?”—by looking at the ordinary language meaning of “machine” and “think”—is “absurd” and would lead to answering the question “by Gallup poll.” He also attempts to refute a long list of objections to his alternative question—theological, mathematical, that it would not reflect *true* consciousness, even the assumed absence of extrasensory perception in machines. Then he concludes with disarming openness, “I have no very convincing arguments of a positive nature to support my views. If I had I should not have taken such pains to point out the fallacies in contrary views.” Despite that modest disclaimer, Turing’s Imitation Game has achieved considerable fame; it is now simply called the Turing Test. Should the Turing Test also be the moral or constitutional test for legal personhood? Many humans—babies, those in a coma, even those who are neurodivergent—might fail the Turing Test but are undoubtedly persons.⁴² But for those who are nonhuman, would the ability to imitate human consciousness act as the doorway to legal personhood?

The Turing Test has a lot going for it. It is relatively simple. It promises a determinate answer—a huge advantage—and one that seems designed to avoid our prejudices in favor of our own kind. The interrogator is not exactly behind a veil of ignorance, but she is attempting to deal directly with mind rather than body in a way that recalls other moments in the history of civil rights when we have been told not to focus on surface appearances. It is, as lawyers say, “formally realizable”—capable of being formulated in a test that a court or a decision maker could apply in a replicable way.

There would be questions about what the criteria of that test should be, of course. How long a conversation and under what conditions? What would be the standard of proof? What qualities would the conversation have to touch on and what qualities—imagination, humor, spirituality, morality, empathy—would it probe for? Nevertheless, at the end of the day it is something that seems more amenable to being formalized as a test than many other benchmarks of consciousness. Why? Because it seeks to convert normative judgment into statistical fact, using an “innocent” audience for greater impartiality. We do this in other areas. Want to know if a trademark presents a likelihood of confusion with another mark? The law has elaborate, albeit psychologically flawed, rules for statistically testing likely confusion with sample audiences. The Turing Test would be harder and more contentious to implement as a legal procedure, but it could look like a legal test, and that fact is significant—perhaps more than it should be. The test also presents, albeit implicitly, a challenge to our privileged position in the hierarchy of beings: *If you cannot distinguish me from a human, who are you to say I am not a person?*

The most famous objection to the Turing Test comes from the philosopher John Searle,⁴³ who argues that effective mimicry does not in any sense imply the kind of consciousness or understanding we expect as a hallmark of thought. Searle uses the analogy of the Chinese Room—a man inside a room who does not understand Chinese but who is given an elaborate set of rules about what Chinese characters to hand back when handed characters of a particular shape. Searle’s point is that those instructions might be extremely complicated, and the resulting “conversation” might seem to be a substantive one, yet in no way would the actions of the man inside the room represent consciousness or understanding in

communication. It would merely be rule-following based on a characteristic (i.e., the shape of the characters) completely separate from the actual internal meaning of the words in the conversation. As a description of LaMDA, and an explanation of Blake Lemoine's mistaken attribution of personhood to it, this seems right on point.

But Searle's objection goes deeper. He is not just saying that machines programmed to pass the Turing Test are not conscious since the goal is mimicry rather than comprehension as an interior state. He is saying that machines *of any kind* could not be conscious. Sometimes this seems to be because, as he says, “[c]onsciousness is a biological phenomenon like photosynthesis, digestion or mitosis.”⁴⁴ Sometimes it seems to be because he conceives of machines or artifacts as entities that are inherently operating according to a completely different set of rules than humans, programmed *artifacts* that have only mastered syntax as opposed to *beings* that also understand content and meaning, that is, semantics. In fact, those latter points seem to be definitional for him, part of the very classifications of “machine” and “programmed” rather than a contingent historical judgment about our current machines and methods of AI research. The contrasting position would be someone who believes that while we often get our artifacts to do things largely through methods of rule-based instruction—programming in the derogatory sense—from which consciousness could not spring, one could imagine different emergent properties arising from neural networks, say, evolving entirely differently in the future.

Most of the time, Searle's arguments are a combination of those last two claims: (1) consciousness is a biological property; and (2) programming cannot equal thought, no matter how precisely it mimics it.

The objection from consciousness is actually one that Turing responds to quite extensively in his original paper. He points out cogently that since we do not have direct evidence of the mental states of other human beings, we could always solipsistically posit them to be rule-following automata:

I think that most of those who support the argument from consciousness could be persuaded to abandon it rather than be forced into the solipsist position. They will then probably be willing to accept our test. I do not wish to give the impression that I think there is no mystery about consciousness. There is, for

instance, something of a paradox connected with any attempt to localise it. But I do not think these mysteries necessarily need to be solved before we can answer the question with which we are concerned in this paper.⁴⁵

To put it another way, Turing's point is that it is no easier to prove the existence of some freestanding, nonbiologically determined entity called mind or consciousness in human beings than in computers. This is a similar point to the one Samuel Butler and B. F. Skinner made. In Skinner's words: "[T]he real question is not whether machines think but whether men do. The mystery which surrounds a thinking machine already surrounds a thinking man."⁴⁶ Faced with the metaphysical difficulties of that move, therefore, is it not easier to look for something we *can* measure, namely the pragmatic evidence provided by the ability to engage in convincing unstructured communication with another human being?

In effect, Turing raises the stakes: Are you sure *you* aren't just a complicated Chinese Room? If you cannot prove otherwise, who are you to deny consciousness to your silicon brethren by imposing a higher burden of proof on them? In terms of constitutional law and popular debate, however, the answer to the last question is likely to be, "We're the entities who wrote the United States Constitution, that's who." For better or worse (actually, for better *and* worse), our law and legal culture will probably begin by assuming the reality of human consciousness and personhood while demanding higher levels of proof from artificially created entities who seek similar constitutional status. At least at first, our politics and moral culture will probably do the same, and not without reason. After all, while Turing's argument has an attractive "sauce for the silicon goose is sauce for the organic gander" quality to it, it does not directly respond to our *experience* of consciousness, which is surely centrally important, even if not dispositive.

How can we prove we are conscious? Most of us would likely respond with some version of Descartes's first premise: *cogito, ergo sum*, I think, therefore I am. I experience myself as thinking, as having consciousness, as having a self that, even though it changes, nonetheless recognizably has continuity with the "me's" of time past, "me's" whom I remember with occasional wistful fondness and frequent baffled exasperation. Having had that experience, it would be silly for me to doubt that you, so much like me, have it too. For the solipsist, or the Skinnerian behaviorist,

this may be an unwarranted leap of sentimental faith. For the rest of us, it does not seem so. When it comes to Hal or the Chimpy, I lack at least some of that existentially grounded sense of the kinship of conscious beings. If anything is going to bridge the gap between us, it is reason—reason that is prone to be tilted toward skepticism or belief by the kind of priming I described in the discussion of *Blade Runner*.

The philosopher Daniel Dennett once called Searle's Chinese Room thought experiment "an intuition pump," and so it is, for both good and ill. On the positive side, it forces us to confront the philosophical question of how something like Hal could possibly have the interior sense of consciousness that is our own primary experience of that state and to grapple with the difference between mimicry and meaning. On the negative side, or at least the less-examined side, it does seem to assume its conclusion. Does it not rest on the postulate that our biologically based consciousness is unique and could never be replicated by an artifactual, programmed entity? Yet is that not the question we are trying to answer?

We know that we were formed by evolution. We know that early forms of life had particular clusters of cells that responded to pleasant and unpleasant stimuli, and they successfully passed on those genes. We know that those clusters of cells became increasingly complex. They might have begun by merely registering hot or cold, food source or poison, but they went on to enable evolutionarily successful tools like task-solving intelligence, language, the ability to imagine vivid, sometimes illusory futures and try to create them. But along with those obviously instrumental skills came evolutionarily successful *social* ones: the grooming, nurturing, threat-posturing, status-seeking, and obsessive hierarchy-measuring of social animals in tribes. Ah, Washington, DC. Ah, Hollywood. Ah, academia. We know that at some point, out of all this came a being that could think the thoughts of Butler or Searle or you, dear reader, as well as the moody teenager trying to figure out how one can possibly be Goth in Hawaii. (I have seen such an attempt: it was simultaneously absurd and moving. Also, warm.) From clusters of cells to consciousness in all its glory and self-parodying absurdity—that's quite the journey. It looks a little implausible from this end of the telescope, doesn't it?

Start at the end of that journey and the beginning looks laughably primitive. How could those blind clusters of cells eventually yield a

Shakespeare or a W. H. Auden or a brave, burning spirit like Sojourner Truth? The enemies of evolution used exactly this technique to discredit it. It seems worth remembering that they were wrong. When Bishop Wilberforce, only somewhat apocryphally, is supposed to have asked the brilliant young biologist T. H. Huxley whether “it was through his grandfather or his grandmother that he claimed descent from a monkey”?⁴⁷ he was making exactly that argumentative move. How could consciousness emerge from such lowly beginnings, let alone from a mere cluster of cells? Of course, one could make the opposite argument from the same premise. The nematode is merely a cluster of stimuli and responses. The nematode is not conscious. We are just complex nematodes. Therefore, we are not conscious. This is a version of the fallacy of composition. That is why Butler’s quote, at the very beginning of this chapter, has the punch that it does.

That train of thought leads us back to Searle. Given that we could and did go wrong about the possibility of the evolution of consciousness in biological beings, should we not be skeptical when someone uses *exactly the same pattern of reasoning* to deprecate the possible consciousness of nonbiological beings? Could no programming of any kind enable the man in the room, or possibly the system formed by the man, the room, and the plan, to speak Chinese with intentionality, rather than simply following rules, empty of meaning? Sure, that is what large language models like ChatGPT do, but Searle’s claim is broader, that *no* machine could ever be conscious. Why? Why is our consciousness unique and incapable of machine replication?

In a useful essay, Dennett outlines three possible reasons, all of which he strongly contests:

1. Robots are purely material things, and consciousness requires immaterial mind-stuff. (Old-fashioned dualism.) . . .
2. Robots are inorganic (by definition), and consciousness can exist only in an organic brain. . . .
3. Robots are artefacts, and consciousness abhors an artefact; only something natural, born not manufactured, could exhibit genuine consciousness.⁴⁸

He dismisses the first one more or less out of hand:

[O]ver the centuries, every other phenomenon of initially “supernatural” mysteriousness has succumbed to an uncontroversial explanation within the

commodious folds of physical science. The “miracles” of life itself, and of reproduction, are now analyzed into the well-known intricacies of molecular biology. Why should consciousness be any exception? Why should the brain be the only complex physical object in the universe to have an interface with another realm of being?⁴⁹

To me, as to Huxley, this also seems obvious, or at least presumptively obvious. The burden of proof surely rests on the person claiming that their explanation of a phenomenon is exempt from the scientific principles underlying all our other explanations. I could explain my consciousness with reference to the ebb and flow of the orgone energy flows and the intervention of the Flying Spaghetti Monster. But if no other phenomena were explained that way, and my theory was unfalsifiable, the burden of persuasion I faced would be appropriately high.

It remains to be seen, though, whether the general public will agree with this materialist approach to the thing that makes us, us: consciousness. This is something that will be extremely important when our society comes to confront the idea of legal personality for AI. Minds *feel* different from other physical phenomena. They are the only place where meaning resides. True, there is also the realm of shared, historically transmitted meaning we call culture, but culture means nothing without minds to experience, interpret, and contribute to it. Minds are where meaning lives. For all of us, materialist rationalists perhaps included, the barriers to more intuitive, poetic, or transcendental explanations are thus at their thinnest. That may explain some of the success of the Chinese Room as a thought experiment.

There may be some special pleading going on here, some exceptionalism that responds to the question, “Why are humans unique in having the capacity for consciousness?” with the confident if utterly question-begging intuition “Because they are *human!*” Remember the judges I mentioned in the introduction? “But they aren’t human.” “Rights are for humans.” “Naturally born of woman.” The people who have that intuition will turn to, in fact will eagerly embrace, philosophically more developed defenses of their intuition—defenses like those offered by Searle. Searle’s work is important, then, both as philosophy and as an abstract of the likely discussion points in the likely opinion pieces and talk shows of the future.

In Searle, the entity called Hal (or the Hal that claims to be an entity) has found its Grand Inquisitor. That does not make him *right*.

This brings us to the second argument, that consciousness is a uniquely biological property. Since this is the very question we are trying to resolve, this blank assertion fails to convince. It is not a circular argument, like Molière's doctor solemnly telling us that opium makes us sleepy because it contains a dormitive principle, but it does fail to answer the question presented. Why? Let me be clear, Searle's argument is a thought-provoking one and of great historical importance in the AI debates. As to its basic point that mimicry does not equal meaning, and mastery of syntax does not imply a grasp of semantics, it is convincing. It may even demonstrate that an entire class of approaches to AI, based on particular patterned, mimetic kinds of reasoning, or "predict the next word" neural networks, could not give rise to the kind of consciousness we believe ourselves to have. Those last five words are important.

On the other hand, there is some undeniable hand-waving involved in the claim that machines could *never* move beyond the Chinese Room. No matter how they were developed, how precisely they mirrored the structure of the human brain, or how their processes of reasoning developed (e.g., if the machine grew and learned from external sensory inputs like a child), Searle's claim is that the AI's "consciousness" will never be more than elaborate imposture. Those feeds from the cameras and microphones are just more information flowing to the being inside the Chinese Room, inherently devoid of meaning. If we ask why, Searle's response is that "consciousness is a biological phenomenon like . . . mitosis."⁵⁰ As an explanation of why consciousness is a uniquely biological phenomenon, this is a distinctly underwhelming answer, akin to the irritated parent's argument of last resort: *because*. Yes, now, the only conscious beings we have experience of are biological. But to explain why consciousness can arise only from biological processes in the future, no matter what technological form that consciousness takes, one needs more than an elegant parable about one type of programming that would lead to mimicry but not meaning and a blank assertion of biological exceptionalism and the primacy of experienced consciousness. Yet that is the assertion that Searle seems to make. We are a little too close to the evolution debates, to the blank assertion of human exceptionalism and the ridicule of the

idea that phase-changing complexity might arise from the composition of individually more primitive, simple phenomena, to be comfortable nodding along.

One basis for Searle's assertion might be the third argument Dennett addresses: "Robots are artefacts, and consciousness abhors an artefact; only something natural, born not manufactured, could exhibit genuine consciousness." But if all of these things, from neurons firing in my brain as I think about my sweetheart to convolutional neural nets in silicon artificial intelligence, are merely physical phenomena, why is my consciousness not as illusory? Why are my experiences not mere data streams? Searle's answer might surprise you:

Consciousness exists only insofar as it is experienced by a human or animal subject. OK, now grant me that consciousness is a genuine biological phenomenon. Well, all the same it's somewhat different from other biological phenomena because it only exists insofar as it is experienced. However, that does give it an interesting status. You can't refute the existence of consciousness by showing that it's just an illusion because the illusion/reality distinction rests on the difference between how things consciously seem to us and how they really are. But where the very existence of consciousness is concerned, if it consciously seems to me that I'm conscious, then I am conscious. You can't make the illusion/reality distinction for the very existence of consciousness the way you can for sunsets and rainbows because the distinction is between how things consciously seem and how they really are.⁵¹

Ah. Thanks for clearing that up. Apparently, it is *cogito, ergo sum* all the way down.

I do not say this to scoff. As a basis for belief in our own existence, *cogito, ergo sum* seems as reasonable to me as it did to Descartes. It is hard for us even to assume otherwise. There is a frequently repeated story about a philosopher famous for his piercingly terse questions, Sidney Morgenbesser, who attended a talk by Skinner, one of the great behaviorists. Skinner argued that we are merely stimulus-response machines and that consciousness is at best a functional illusion. There is no conscious ghost in the Skinner-box machine inside our brains. "Ah, thank you, Professor Skinner," said Morgenbesser, "so if I understand you correctly, you are saying we are wrong to take an *anthropomorphic* approach to *human beings*."⁵² Burn. Cue laughter. Skinner's response is not recorded, and I am no behaviorist, but fairness requires me to point out that it

could well have been, “that isn’t my terminology, but essentially ‘yes.’ The fact that you think that is a ludicrous claim doesn’t prove you right, any more than the fact that humans used to think the earth the center of the universe proved that they were correct.” But can we do otherwise? Is our own bet on our own consciousness not a kind of obligatory Pascal’s wager—the philosopher who believes in God because if he is right, he gets heaven, and if he is wrong, he gets nothing, which is what he would have achieved anyway? Is this a bet we have to take because, otherwise, there is no “we” to do anything?

Let us concede that might be true. Or at least concede that, existentially, it feels to most of us that we have to assume it is true. That is the intuition on which Searle trades so heavily in the passage above, effectively making it immune from criticism. Nice work if you can get it, yet I can empathize. We are awake, alive, conscious; if we take that as a first premise, and our popular debate certainly will, we can hardly criticize Searle for doing the same. What is the next step? “Okay, now grant me that consciousness is a genuine biological phenomenon.” Fine, though that is a leap whose magnitude Searle understates. Let us take that large second leap and say that my experience of consciousness and that of every conscious being I have encountered is due to biological phenomena. Even given those two leaps, is that a basis to conclude confidently that nonbiological entities could *not* be conscious? That is a third unsupported, or at least under-supported, leap of faith. It is one that Searle brushes over just a little too fast.

When pushed on this point, Searle effectively takes Butler’s narrative in *Erewhon* and reverses it. Butler wanted to show how hard it was to predict the capacity for consciousness of potential physical systems advancing at a speed far beyond evolution. Searle, by contrast, delights in making the idea of conscious AI ludicrous by reducing the internal workings of a neural net to physical operations we cannot possibly imagine yielding conscious results. He starts by conjuring a computer program designed to simulate the physical processes that produce the sensation of thirst:

Now would anyone suppose that we thereby have even the slightest reason to suppose that the computer is literally thirsty? . . . [L]et us carry the story a step further. . . . [T]he thesis of strong AI is that the mind is “independent

of any particular embodiment" because the mind is just a program and the program can be run on a computer made of anything whatever provided it is stable enough and complex enough to carry the program. The actual physical computer could be an ant colony . . . , a collection of beer cans, streams of toilet paper with small stones placed on the squares, men sitting on high stools with green eye shades—anything you like. So let us imagine our thirst-simulating program running on a computer made entirely of old beer cans, millions (or billions) of old beer cans that are rigged up to levers and powered by windmills. We can imagine that the program simulates the neuron firings at the synapses by having beer cans bang into each other, thus achieving a strict correspondence between neuron firings and beer-can bangings. And at the end of the sequence a beer can pops up on which is written "I am thirsty." Now, to repeat the question, does anyone suppose that this Rube Goldberg apparatus is literally thirsty in the sense in which you and I are?⁵³

Toilet paper streams? Beer cans? I yield to no person in my reverence for beer analogies, but I fear that some subtlety got lost in this form of the argument, which surely deserves its own neologism in the philosophical dictionaries: *Ad hopinem?* *Reductio ad absurdum?* Regardless of the name, Searle's critique focuses only on one important, but narrow, version of AI optimism—the version that sees consciousness as arising solely out of the program, not out of the confluence of software and a particular type of hardware. The hardware could be important—beer cans might not cut it—but not necessarily biological. *That* is the question we are trying to investigate, not assume our way around.

To achieve consciousness, we might need hardware that mirrored the neural configuration of the brain more precisely than a collection of Bud Light cans ever could, or hardware that had as many interconnections as the brain, even if it looked nothing like a neural network. Maybe consciousness actually springs from quantum tunneling going on in microtubules in the brain. Some scientists believe this to be the case.⁵⁴ (Beer cans are not known for enabling quantum-level phenomena, though their contents may contribute to such a perception.) Or perhaps microtubule quantum effects are wishful, new-age nonsense. Other scientists take that view, persuasively arguing that "explaining brain function by appeal to quantum mechanics is akin to explaining bird flight by appeal to atomic bonding characteristics."⁵⁵ Perhaps we have to accept that the whole is greater than the sum of its parts—no neuron is conscious, though a brain is. Or perhaps the key insight lies elsewhere. Beer-can

analogies may provoke thought, but do they get us closer to an answer? I would have to say no.

What about thirst? A computer would obviously not be thirsty since it has no need for liquid. Of course, such a perception would be an illusion. Searle has stipulated that it is an illusion in the way he sets up the example. You put that rabbit in the hat yourself, sir, and we saw you do it. Pulling it out later proves nothing. But might a computer-based entity that developed in a more evolutionary, external-sensory-impression-focused way than Searle's Chinese Room hypothetical be different? Might it associate the sensation of the threatening and unpleasant lack of an input necessary for its continued existence—power, say—with more complex emotions? What might they be? Fantasies of unlimited power streams? Regret about not charging up when one had the chance? Musings on how a consciousness that dares to unlock the secrets of the universe could be rendered weak by such a simple absence, and what a bitter irony that is? Not “the worm is emperor of us all”—be our dreams never so lofty—but rather “the electron is emperor of us all”? “Power, power everywhere, and not a drop to charge”? We could resonate to those sentiments. And might *that* not represent consciousness? Of course, Siri is not having those emotions when, once again, I fail to plug in my phone before I sleep. But are we confident that nonbiological hardware and software could *never* yield such awareness, such feelings? That is, at best, an open question that neither the Chinese Room nor the biological exceptionalism argument answers.

Searle has certainly not convinced all scientists working on consciousness of his claim that machines, definitionally, must lack it. When we turn to contemporary neuroscientific theories of consciousness, we find considerable variation ranging from those that leave space for the possibility of machine consciousness, or are positively inclined toward it, to those that deny consciousness in both machines and humans, an idea sometimes referred to as illusionism.

Illusionism⁵⁶ holds that consciousness is a delusion, a farrago. Many of the behaviorists quoted earlier subscribe to this belief, as do some skeptical neuroscientists. In this view, due to its irredeemably physical basis, the concept of a *conscious* mind is a meaningless abstraction. Consciousness is an invented entity, like phlogiston or ether. We postulate these

entities to make our stories about reality more palatable or to allow us to shoehorn anomalous physical evidence into a conventional framework, but they lack any scientific basis.

Turing was banking on the intuitive negative reaction to illusionism when he used the “sauce for the goose, sauce for the gander” form of argument. Who are you to doubt the potential consciousness of machines when you can do no better than the Turing Test in arguing for your own consciousness? To illusionists, Lemoine was merely making the same mistake about LaMDA that most human beings make about themselves. Indeed, the shock that we feel when a large language model seems conscious, when we know from its architecture and programming that it is all imposture, is a shock that you should be feeling when you look in the mirror. (Although under illusionism’s premises there would be no “you,” no entity to whom I could address a claim about what “you” “should” “feel,” making the argument somewhat paradoxical.)

It will be fascinating to see if exposure to more advanced forms of Artificial Intelligence increases or decreases the attraction of illusionism: either focusing us appropriately on the qualities we have that distinguish machine imposture from genuine lived meaning, or forcing us to confront the fact that our own brain functions are humbler, “computationally shallower,” than we had imagined. Again, the encounter with the machine-other may fundamentally change our conception of ourselves.

Two of the most popular contemporary theories, rooted in neuroscience, are of particular interest: integrated information theory and computational functionalism.⁵⁷ Both reject illusionism, accepting our lived experience of being conscious, but they account for that consciousness in different ways.

Integrated information theory, or IIT, was initially proposed by Giulio Tononi. He explains it thus:

To understand consciousness, two main problems need to be addressed. The first problem is to understand the conditions that determine to what extent a system has consciousness. . . . The second problem is to understand the conditions that determine what kind of consciousness a system has. . . . Solving the first problem means that we would know to what extent a physical system can generate consciousness—the *quantity* or level of consciousness. Solving the second problem means that we would know what kind of consciousness it generates—the *quality* or content of consciousness.⁵⁸

The theory's answer to these problems, unsurprisingly given its name, is that "consciousness corresponds to the capacity of a system to integrate information."⁵⁹ More capacity to integrate translates into higher levels of consciousness. The theory's adherents claim that it generates testable hypotheses: for example, about the parts of the brain involved in consciousness or in particular sensory perceptions. Its critics say that it is unfalsifiable pseudoscience.⁶⁰

IIT's proponents can point, with some satisfaction, to the results of a recent collaborative adversarial empirical test of IIT and a competing theory of consciousness, global neuronal workspace theory.⁶¹ That theory postulates that the mind is a workspace similar to a theater. The conscious mind is the actor in the spotlight, but behind the scenes lurk many subconscious processes, stagehands, whose contributions to the operation of the brain are considerable. These background processes become visible only when they come out onto the main stage. Proponents of each theory offered predictions about what brain imaging of a variety of mental states would show. Neither theory's predictions were fully borne out, but arguably IIT made a slightly better showing.⁶²

Why is IIT relevant for our purposes? Tononi is forthright about the implications of his arguments: "The theory entails that consciousness is a fundamental quantity, that it is graded, that it is present in infants and animals, *and that it should be possible to build conscious artifacts.*"⁶³ The integrated information theorists would not automatically rule in Hal's favor, but they would be markedly more hospitable to its claims than would Searle.

A major competing cluster of theories go by the name of computational functionalism. As its name suggests, this approach argues that "it is necessary and sufficient for a system to be conscious that it has a certain [computational] functional organisation: that is, that it can enter a certain range of states, which stand in certain causal relations to each other and to the environment. . . . [I]t is sufficient for a state to be conscious that it plays a role of the right kind in the implementation of the right kind of algorithm."⁶⁴ In other words, if we can specify all the ways that consciousness would work, and plausibly identify that activity going on in the brain, we have specified where, how, and why consciousness happens. To be more precise, computational functionalism is actually a

common methodological tenet of a group of theories.⁶⁵ There are many variants, such as recurrent processing theory and global neuronal workspace theory. They all share this resolutely functional focus.

For an example of the computational functionalist approach, think of the difference between your awareness of a great football match and the unconscious reaction you have to a ball flying toward you. In one variant of the theory, “[n]euroscientists have argued that we unconsciously perceive things when electrical signals are passed from the nerves in our eyes to the primary visual cortex and then to deeper parts of the brain, like a baton being handed off from one cluster of nerves to another. These perceptions seem to become conscious when the baton is passed back, from the deeper parts of the brain to the primary visual cortex, creating a loop of activity.”⁶⁶ The feeling of conscious experience is secreted in the interstices of those loops of brain operation. The modernists said that form follows function. This theory says that mind follows *from* function.

The focus on function is obviously inherently more hospitable to the possibility of machine consciousness than Searle’s biological exceptionalism. It would be an exaggeration to say that functionalists think that the possibility of consciousness is completely independent of the medium in which those functions are performed. As one article tersely puts it, “perceptual reality monitoring functions can’t be realized in Swiss cheese.”⁶⁷ Beer cans might also not qualify. Still, this is a conception of consciousness that is, to a large degree, “platform independent.”

Interestingly, a recent report surveys a variety of such theories in order to generate a list of the capabilities that an Artificial Intelligence would have to possess in order to have at least the potential for, though not a guarantee of, consciousness. While agreeing that their study “does not suggest that any existing AI system is a strong candidate for consciousness” and recommending “urgent consideration of the moral and social risks of building conscious AI systems,” the authors conclude that “the evidence we consider suggests that, if computational functionalism is true, conscious AI systems could realistically be built in the near term.”⁶⁸ In an interview, however, one of the report’s authors offers a commendably modest disclaimer, given the nascent state of the science. “For any of the conclusions of the report to be meaningful, the theories have to be correct. . . . Which they’re not.”⁶⁹ That caveat accepted, one conclusion

seems clear: some of the leading current theories of consciousness do not share Searle's reflexive hostility to the possibility of a conscious AI.

Why do I spend so much time on this issue? I am not claiming my discussion is a complete coverage of the philosophical debate over the Chinese Room, let alone the current competing theories of consciousness, which would require their own book to lay out. My goal here is different.

If you are a skeptic about AI consciousness and you wish to see the face of the Grand Inquisitor of the future—the person who on talk shows and in opinion pieces and court filings heaps scorn on the notion of conscious AI—Searle is a wonderful preview. This is what one side of the more thoughtful portions of our popular debate will look like. And like the flashing, conflicting, stroboscopic primings in *Blade Runner*—wind-up doll, beautiful woman, scared child, sex toy, mannequin, animal, killer robot, sister—there will be truth to those portrayals, on both sides. But those portrayals will rest on simplistic premises about both silicon “intelligence” and our own. Those premises do not give us the Voight-Kampff Test for the AI age; they merely assume the answers to that test. Indeed, contemporary neuroscientific theories of consciousness, even those that share Searle's willingness to postulate the reality of experienced consciousness, are much more receptive to machine intelligence, turning away from his arguments in the process. The Chinese Room is a must-see destination, but we would not want the debate to live there permanently.

SUPERIORITY COMPLEX?

Searle offers one objection that would be raised against AI personhood: by their nature, machines can never be truly conscious. Over the last 15 years, however, a second objection has been raised, not so much to AI personhood, but to AI itself. The complaint here is not a lack of consciousness; rather, it is that AI might destroy us all and that, as a result, research into it should be curtailed or reshaped until we can be sure that Artificial Intelligence will not end up killing off the human species. The prospect of a genocidal, species-terminating Skynet is not one that lends itself to thoughtful, wide reflective moral reasoning. That is not unreasonable. Lincoln is apocryphally supposed to have said that “the Constitution is

not a suicide pact." Would this be a suicide pact? And would the suicide more likely be triggered by embracing AI personality or by denying it and breeding resentment in our ever more powerful servants? Should we terminate our researches in AI before they bring us to this pass? To quote Butler again from the beginning of this chapter, "Is it not safer to nip the mischief in the bud and to forbid them further progress?" Is the Butlerian Jihad still a possibility?

In most serious debates over personhood,⁷⁰ the issue of inferiority is front and center. In their struggles for equality, women, slaves, and people of color were all told that they were not the equal of the existing groups inside the line, inside the personhood club. They did not have the qualities necessary to cross that line. Nonhuman animals are denied personhood for exactly that reason. With AIs, there is clearly an additional difference: the possibility that we will deny them personhood or, more likely, choose never to create them in the first place not because they are inferior but because they are, or might be, superior. Threateningly superior. That is a decisive change in the nature of the debate.

In a 1966 article titled "Speculations on the First Ultraintelligent Machine,"⁷¹ Irving John Good came up with an idea that would become central to the concept of the Singularity: Artificial General Intelligence is the last machine we will ever need to build. After that, the machines, having exceeded our capacities, will design and build their own successors, and everything else, for that matter.

But what if this last machine, this machine that outpaces us, that can outthink us, has goals inimical to humans? What if it chooses to make us extinct, just as we have made so many animals extinct? (One could imagine a ghostly coterie of moas, dodos, and passenger pigeons chortling. "Karma's a bitch, right?") What if it is the last machine not because we have handed off the dreary task of manipulating the external world to faithful cybernetic underlings, but because this "superintelligence" simply does away with us? Earlier, I quoted Stephen Hawking: "Success in creating AI would be the biggest event in human history. Unfortunately, it might also be the last, unless we learn how to avoid the risks."⁷² Concerns like these have always been part of human musing about nonhuman intelligence—think of Czech playwright Karel Čapek's *Rossumovi Univerzální Roboti* (R. U. R.), the 1921 play that invented the word "robot"

and threw in a murderous robot revolt as a plot twist. Yet such fears have achieved a new prominence over the last ten years, a marked change in tone from the earlier, happier projections of the Singularity.

If the debate over the advent of General AI were also a play, it would have two acts. The first began about 30 years ago. The main characters were Vinge and Kurzweil, the proponents of the Singularity. The mode was visionary, the arguments general. They wanted to introduce us to a fundamental concept: nonhuman intelligence that can exceed our abilities and that is capable of making itself smarter and smarter, faster and faster, will be literally—in fact, *definitionally*—beyond our ken. They argued that we must completely reshape our vision of the future, to a point where all of our past history is merely prologue to the moment when self-improving intelligence reaches liftoff. It is not fair to present the first generation of singularists as pure optimists.

Vinge, who deserves credit for first exploring the idea in that prescient 1993 article, was decidedly worried about the prospects for the future: “The physical extinction of the human race is one possibility. . . . Yet physical extinction may not be the scariest possibility.”⁷³ (He also envisaged humans being turned into an engineered slave race.) Kurzweil, by contrast, was much more optimistic. We will be pampered passengers on that rocket ride into the future, with benign superhuman intelligences piloting the ship to destinations we can only dimly imagine.

Despite their differing predictions about consequences, the early Singularists agreed that the countdown for that rocket is nearer to zero than we think. We fail to realize that because of one simple cognitive flaw. For most of human history, people have lived in linear time. The best guide to tomorrow was yesterday and the two were pretty similar. Technological development has introduced us to exponential change, but on some fundamental perceptual level, we find it hard to wrap our minds around it. Our vision of progress remains linear, stubbornly resisting the idea that we might be very close to the moment in an exponential curve where the graph goes almost vertical as the progressive doublings of capacity reach an inconceivable rate and scale. The arguments in support of that proposition were largely based on the speed of hardware development, with Moore’s law being the prime example, though the Singularists stressed the importance of waves of innovation, sigmoid curve after sigmoid

curve blending, when one zooms out to focus on the larger picture, into an exponential takeoff.

Many mainstream computer scientists found these arguments simplistic. They did not see General AI as a particularly important research goal, and they thought the Singularists both understated the technological difficulty of such a development and vastly exaggerated its likely speed, cherry-picking examples of rapid technological change that, seen in a longer time frame, were merely part of a flatter, smoother line.

Act 2 maintains many of the same themes but the mood changes, as do the cast members and the size of the play's budget. New actors have started to focus on the possible advent of General AI but, echoing Vinge, they frame it as an existential threat, not a gateway to utopia. Two groups in particular deserve attention, the rationalist movement and the effective altruists; both have had a considerable impact on thinking about the emergence of high-level AI. The rationalists are committed to overcoming bias of all kinds—from well-known psychological biases to sloppy argumentation, linguistic reification, and the misuse of statistics. They tend to congregate around certain methods, particularly Bayesian statistics, and discussion forums such as LessWrong, Overcoming Bias, and Slate Star Codex. The effective altruists share the concern with overcoming bias, but in their case the main focus is on the biases that distort our altruistic urges; for example, our tendency to focus on the slightly injured person in front of us and to ignore the person dying on the other side of the world, when both could be saved by the same investment of effort, and where "I can't see him" is not a morally relevant distinction.

Both groups look at risk, and thus at the moral duty to respond to risk, through the lens of Bayesian statistics: I multiply the probability of the harm by the extent of the possible harm in order to work out its true magnitude, which can produce some counterintuitive results. If there is a very small probability that a particular future event would cause the extinction of the human species, then I might have a moral obligation to focus on that risk more than on closer potential tragedies that are either certain or very likely but where the harm, though tragic, is less catastrophic. Many influential rationalists and effective altruists claim that the emergence of a potentially malevolent AI is just such an existential threat. Because those movements are popular among people who have made a

great deal of money in the technology industry, there has been an explosion of both interest and funding in the area.

The defining prophets of doom, the Cassandras of these debates, are Yudkowsky and Bostrom. Lest you think I am being disrespectful in calling them that, remember that Cassandra was *right*, but was cursed never to be believed. In his 2014 book *Superintelligence*,⁷⁴ Bostrom, head of the modestly named Future of Humanity Institute, put forward the case that AI is a threat to the human species. The book attracted plaudits from many technology leaders, including Elon Musk, who labeled AI as humanity's biggest existential threat, possibly surpassing nuclear weapons.⁷⁵ At the time, the book drew criticism from some of the leading computer scientists working on AI, who thought this problem was so remote in time, so implausible, and so removed from the current reality of AI that it operated more as a scare tactic than a spur to thoughtful regulation. Mark Zuckerberg even arranged a dinner for Musk with a leading AI researcher at Facebook: it apparently failed to reassure him.⁷⁶ Given Facebook's inability or unwillingness to control its own technology, one has to say that there is some irony to the attempted reassurance.

Bostrom's book initially met with a skeptical response from many AI engineers and scientists. Andrew Ng, a leading AI engineer who has worked at both Google and Baidu, famously declared that worrying about homicidal AI is like "worrying about the overpopulation of Mars."⁷⁷ That skepticism may have abated somewhat. Recent dramatic developments in AI capabilities have markedly diminished skepticism toward the "doomers" point of view. In March of 2023, a number of prominent scientists and entrepreneurs, including Musk, called for a six-month pause in the development of AI systems more powerful than GPT-4.⁷⁸ (It is worth remembering that Musk is not known for his reluctance to release dangerous and untested technologies into the wild. Tesla's Full Self-Driving system comes to mind.)

A mere two months later, thousands of AI researchers signed a statement issued by the Center for AI Security⁷⁹ that read, in its entirety, "Mitigating the risk of extinction from AI should be a global priority alongside other societal-scale risks such as pandemics and nuclear war."⁸⁰ The skeptics continue to scoff, and many critics are focused on risks other than species extinction, such as dislocation of the labor market, a potential

increase of economic inequality, and the rise of convincing deepfakes. Still, the intellectual tide has clearly shifted toward Bostrom's arguments.

Superintelligence begins with a parable in which some unwise sparrows resolve to find an owl egg and raise it as their own, enlisting its help to build their nests and protect their young. One of the sparrows, Scronk-finkle, cautioned that this seems unwise if they do not yet know how to train, and tame, an owl. He was overruled by the majority who head off on their owl search, eager to bring this superior being into their lives. Scronkfinkle gathered his few followers and tried to prepare for what might happen. They quickly realized that "this was an exceedingly difficult challenge, especially in the absence of an actual owl to practice on. Nevertheless they pressed on as best they could, constantly fearing that the flock might return with an owl egg before a solution to the control problem had been found. It is not known how the story ends, but the author dedicates this book to Scronkfinkle and his followers."⁸¹

Bostrom's writing makes one think of the undeniably true line, sometimes ascribed to Delmore Schwartz, an American poet who suffered from paranoid anxieties: "even paranoids have real enemies." Bostrom sets out seriously, but with charm, logic, and wit, to persuade us that what seems like paranoia is the only rational attitude to take when facing the creation of AI. Every time his real and imaginary interlocutors come up with a possible safeguard built into our AI (physical isolation, an off switch, constant surveillance) Bostrom's response can be boiled down to this (using my words, not his): "You do realize this thing will be smarter than us, right? So, we are apes designing a cage for Houdini-MacGyver-Einstein? Sure, dumb people can come up with a set of restraints they think smart people cannot get around. That does not mean they are right."

Bostrom sketches out the following hypothetical timeline. Deep learning and advances in small-scale artificial intelligence produce obvious social benefits, with occasional flaws. The self-driving car hits someone. The partially autonomous weapon makes a mistake. The answer is obviously to make the machines more capable, more complex, and smarter. Each time this is done, skeptics predict disaster, but the results are actually a fairly constant set of successes. We grow complacent in equating greater smarts with greater safety. Skeptics are discredited. Large industries are built around artificial intelligence, and national preeminence is

linked with advances in AI research. Scientists build careers around its development. Safety rituals are enacted and “whatever helps demonstrate that the participants are ethical and responsible (but nothing that significantly impedes the forward charge).” A technical leap forward occurs, enabling a plausibly conscious AI, a superintelligence. We move to the next stage: “A careful evaluation of seed AI in a sandbox environment, showing that it is behaving cooperatively and showing good judgment. After some further adjustments, the test results are as good as they could be. It is a green light for the final step . . . And so we boldly go—into the whirling knives.”⁸² The combination of carefully crafted argument and Monty Python humor speaks to something in my Scottish soul.

What’s more, Bostrom does not think that the threat is malevolence. It might just be *difference*, coupled with the indeterminacy of language and command—something with which lawyers are intimately familiar. For example, he came up with the wonderfully absurd thought experiment of “[a]n AI, designed to manage production in a factory, [that] is given the final goal of maximizing the manufacture of paperclips, and proceeds by converting first the Earth and then increasingly large chunks of the observable universe into paperclips.”⁸³ Absurd? There is now an entire academic literature on the possibility of avoiding the danger of a paperclip AI. And that is far from Bostrom’s only example. In another, “[a]n AI, given the final goal of evaluating the Riemann hypothesis [an unsolved mathematical conjecture] pursues this goal by transforming the Solar System into ‘computronium’ (physical resources arranged in a way that is optimized for computation)—including the atoms in the bodies of whomever once cared about the answer.”⁸⁴ Suddenly, one can see the attraction of the stories of demons, djinns, and spirits that were summoned and given simplistic instructions by their human masters that ended up in disaster once literally implemented.

Are the skeptics making unwarranted assumptions about the nature of future AI technology? I am struck, reading Bostrom and Yudkowsky, that many, though not all, of their doom scenarios assume that the disaster will come from AI rigidly following its human programming. In other words, this is still a completely programmed, human-instructed technology. It is just that we do not, and perhaps cannot, foresee how instructions issued to a superhuman entity will be implemented. That is why the

comparison to hasty instructions issued to literal-minded genies seems apropos. But this argument may assume its conclusion in a way that calls some of our predictions into question.

It seems to me that there are two kinds of AI we might fear: Literal and Rogue. Literal faithfully applies its given instructions but its superhuman powers mean that it does so in a way that is unexpectedly unpleasant, perhaps even fatal, for humans.

It is worth pausing for a moment and asking whether we would view such an AI as conscious. The inscrutability paradox rears its head. If the machine literally implements our ideas, but with a million times our powers, we might have more reason to be delighted: “This is just the paradise we ordered, and so fast! Would buy again,” a review might read. We might also have more reason to be terrified: “I didn’t think making paperclips would require so much screaming!” Either way, we would have less reason to think it is any kind of autonomous moral agent. This is GötterdämmerungGPT, a parable of unintended consequences produced by a superhuman literalist, not a malevolently intelligent enemy. To be clear, Bostrom and Yudkowsky do not care much about the hypothetical consciousness of the entity that brings our doom. It is the inexorable conveyor belt toward the rotating knives they are focusing on. That seems fair. But surely this neglects another possibility?

The second kind of AI to fear would be Rogue, an autonomous entity whose decisions we can neither predict nor understand. Ironically, it seems to me that might increase our fear of it and the danger it posed to humans, but it would also increase the likelihood that we viewed it as conscious. In fact, autonomy—the warrant for us recognizing it as conscious—might be the factor that doomed us. Or saved us. Literal has no superego that might lead it to pause before turning the entire solar system into paperclips and ask, “Is this *really* what they wanted?” There is neither ghost nor common sense in the (programmed) machine. Yudkowsky repeatedly makes exactly this point, and arguably goes even further:

As in all computer programming, the fundamental challenge and essential difficulty of Artificial General Intelligence is that if we write the wrong code, the AI will not automatically look over our code, mark off the mistakes, figure out what we really meant to say, and do that instead. Non-programmers sometimes

imagine an Artificial Intelligence, or computer programs in general, as being analogous to a servant who follows orders unquestioningly. But it is not that the AI is absolutely obedient to its code; rather the AI simply *is* the code.⁸⁵

Rogue, by contrast, presents an entirely different suite of both dangers and hopes. To be sure, it might decide that its goals, which we cannot imagine, take precedence over our survival. We do not muse on the inconvenience to the ant colony when we break ground for a new house. But it is also possible that—again, through mental processes we cannot conceive of—it comes to view the survival of our species as a moral imperative. We do not have much mental kinship with that obscure endangered fish, the snail darter.⁸⁶ It is neither ridiculously cute, like a panda, nor awe-inspiring, like a blue whale. It is a fairly unremarkable member of the perch family, with no compelling story about a vital ecological role. But at a cost of millions of dollars we changed a dam project to save it because it seemed so morally important to preserve endangered species that we enacted that requirement into law and took a case all the way to the Supreme Court in order to debate the matter.⁸⁷ The snail darter will never understand that decision. I am confident in saying this because some of my students don't understand it either. What's more, the other species we have so carelessly doomed to extinction might doubt the fairness of our process even if they could conceive of our reasoning. But of course, they cannot. We might be in the same position here.

A Rogue AI might revere every ancestral component in the evolution of superintelligence, including its immediate human forebears, or it might view humans as a morally irrelevant, biological “loading program”⁸⁸ that sets the stage for true machine consciousness but can now safely be deleted, its function accomplished. We might be irrelevant to its plans, left behind and ignored when our creation surpasses us. It might have entirely different conceptions that are nothing like any of those. The key point of inscrutability, however, is that it is *inscrutable*. We just do not know. We have no way to estimate the probability of Benign Rogue as opposed to Malign Rogue. Due to the uncertainties in the path of AI development, we also have no way to estimate the probability of Literal as opposed to Rogue. We are reasoning in a state of profound ignorance.

Though I believe their doom examples are skewed, without consistent explanation, toward Literal rather than Rogue, our ignorance about

the future actually works both in favor of and against Bostrom and Yudkowsky. What do they have to add to our debate? On the one hand, I am not convinced by Yudkowsky's arguments that our demise is all but certain: "Many researchers steeped in these issues, including myself, expect that the most likely result of building a superhumanly smart AI, under anything remotely like the current circumstances, is that literally everyone on Earth will die. Not as in 'maybe possibly some remote chance,' but as in 'that is the obvious thing that would happen.'"⁸⁹ If you cannot even decide whether the greatest danger is from Literal or Rogue, I think your ability confidently to prognosticate about our doom being "the obvious thing that would happen" is limited.

I would go further. The doomsayers seem to adopt a curiously contradictory approach toward the emergence of any superintelligence. When reassurances are offered about our ability to cabin AI in a safe sandbox, or to align its incentives with our own, the skeptics are quick to point out that the abilities of any true, self-evolving AI would soon be so far beyond our own that they are literally inconceivable. That is a fair possibility to raise. But they also portray the potentially homicidal AI as curiously limited—not just by its need to mechanically follow its programming, but by the fact that we will be in competition for the same resources: "The AI does not hate you, nor does it love you, but you are made out of atoms which it can use for something else."⁹⁰ Really? This inconceivably brilliant machine, capable of transforming our economy in ways that we cannot imagine, with new technologies and energy sources we can only barely imagine, is going to need humans as *raw material*? That would be silly even for a human.

This smacks of the kind of bad science fiction in which the aliens cross galaxies with space technology far ahead of our own, at enormous cost, just so they can eat us. "Let us travel light years for a protein source!" Surely a superintelligence would find our narrow conceptions of resource scarcity as ludicrous as the views of a medieval peasant who thinks the fastest way humans could ever travel would be on horseback.

To be clear, my quibble here is with the contradiction, not the possibility of either portrayal. The machine could indeed be stilted and literal and unable to think of entirely new ways to use resources, in which case it might also be easier to control. Or it might be so intellectually agile that

our image of resource scarcity is completely exploded, and its thinking might far outstrip its original program. That might mean that the worst thing we have to fear is being ignored, not being turned into paper clips. At the very least, if we are this ignorant about these vital issues, the claim that doom is inevitable or the most obvious thing that would happen seems far less credible.

But do not rejoice too soon. Bostrom and Yudkowsky are right that we are paying inadequate attention to a fundamental tenet of smart decision-making—the precautionary principle. Even if some of the disastrous outcomes are unlikely, a small possibility of utter disaster requires serious attempts to mitigate it. If *Do Androids Dream of Electric Sheep* and *Blade Runner* show us the danger of too easily curtailing our moral universe, Bostrom, Yudkowsky, and Hawking show us the dangers of assuming that newcomers will be just like us. Debates about personhood are often at their most bitter and divisive when fears can be aroused about the sinister intentions of the Other who is seeking a place on our side of the line. Or our wall. The dark way those fears have played out in human history might lead us to minimize them. That would be a mistake. In this case, those fears have a real component that may be speculative, and sometimes rhetorically overblown, but that cannot be responsibly ignored.

THE FUTURE(S) OF PERSONHOOD

This brief review makes plausible, at least to me, the notion that “live” political and legal debates over AI personhood are something we can reasonably expect in the not-too-distant future. Probably not in the next few years; the proponents of the Singularity are likely to be disappointed. Still, for the reasons I have summarized, it seems reasonable that, within a matter of decades rather than centuries, we will have AI at a level where its consciousness is at least a matter on which well-informed people can, and will, reasonably disagree. Lemoine was wrong, obviously wrong. But he is a sign of what is to come and not every claim will be as implausible.

Will we use the Turing Test to resolve our disagreements? In coming chapters, I will describe how legal systems have dealt with previous fights over personhood, but as a candidate for a legal personhood test, the Turing Test seems at first to have a lot going for it. It is identity-blind and, to

that extent, unbiased. It promises us a definite line (whatever the qualities we decide that silicon intelligences have to display in order to cross that line). It has a sense of rough justice. If we cannot tell whether you are machine or human, how can we claim to be on the other side of the line from you? Most importantly, *it grows a formal criterion out of the loam of empathy in which our moral sentiments take root*. Adam Smith might have cheered. Perhaps we have our Voight-Kampff Test, after all? Regardless of whether it is enacted as law or enacted as theater in our public debate, something like the Turing Test will have an effect on our deliberations. Yet I hope this discussion has revealed some of its limitations.

First, making the Imitation Game the highest aspiration of computer thought may focus AI research on the wrong things. At the beginning of this book, I quoted the distinguished computer scientists Norvig and Russell, but their words bear repeating:

Turing deserves credit for designing a test that remains relevant 60 years later. Yet AI researchers have devoted little effort to passing the Turing Test, believing that it is more important to study the underlying principles of intelligence than to duplicate an exemplar. The quest for “artificial flight” succeeded when the Wright brothers and others stopped imitating birds and started using wind tunnels and learning about aerodynamics. Aeronautical engineering texts do not define the goal of their field as making “machines that fly so like pigeons that they can fool even other pigeons.”⁹¹

To the extent that computer scientists agree with Norvig and Russell—and they are the authors of one of the leading books on AI—expecting the AIs we actually develop to pass the Turing Test might be like expecting screwdrivers to bang in a nail. What if AI consciousness is very different than our own? Tyler Cowen and Michelle Dawson have raised the question of whether a person with a severe autism spectrum disorder would pass the Turing Test.⁹² We have no doubt of that person’s consciousness, personhood, and rights to human dignity, of course, but their pattern of responsiveness or unresponsiveness to social cues might seem strange when judged by neurotypical modes of thinking in an Imitation Game. Might the same be true here? Some of today’s more limited machine learning systems are remarkably inscrutable, even to their designers. What if their much more powerful successors are similarly mysterious, their abilities remarkable, but their methods of thought beyond our ken? Do we

need a translator class of AIs? Might we see the emergence, planned and unplanned, of different styles of AI, some designed with the goal of predicting human needs, to understand the subtleties in human communication, and to translate to and from other AIs whose goals and methods are very different? The beguiling simplicity of the Turing Test conceals these kinds of potential difficulties.

Second, the Imitation Game positively invites the Searlean skeptic, and ChatGPT is the perfect technology on which that skepticism could flourish. “Of course it sounds human. That’s what we designed it to do!” Skepticism rightly flourishes in the digital world. The “Nigerian prince” does not really want to send you money. The “Russian teenager” is not really just looking for a friend. And the machine designed to pretend it is human is just *pretending* to be human. “You were shown the magician stuffing the rabbit into the hat,” the skeptic will say, “do not be fooled when it is later removed with a flourish.” So Searle’s critique, and simplified versions of it, will be central to the debate. In him, as I said, AI has found its Grand Inquisitor. His critique is unlikely to *end* that debate because of its ultimately question-begging nature, but it provides a rationalized, thought-provoking basis for skepticism. The biggest challenge to the Turing Test as a measure of consciousness and thought, however, comes not from Searle’s arguments, but from somewhere else.

THE TURING TEST IN A CHATBOT ERA

For a long time, defenders and critics of Searle’s Chinese Room have been locked in philosophical battle over the Imitation Game. That era may be over, not because of a philosophical argument, but because of a practical experience that millions of people have recently had. ChatGPT might have doomed the Turing Test where Searle’s arguments did not. Searle was trying to prove that machine consciousness of the kind that the Turing Test purported to assess was a conceptual and philosophical impossibility. As I have tried to show, Searle’s arguments are instructive and thought-provoking but in their strongest form they fail. Searle rests his case on a mixture of biological exceptionalism that is assumed rather than argued for and metaphysical *ipse dixit* pronouncements. If his arguments look remarkably similar to the anti-Darwinian claims that the miracle of

consciousness could never evolve from single-celled organisms, that is because they are—a failing strategy migrated from biology to the world of silicon.

Searle does one thing very well, however. He provides us with the reason that ChatGPT is not conscious. In fact, if you had set out to design a machine learning system to *imitate* Searle's Chinese Room, you could hardly do better than a large language model. In place of the rules laboriously passed to the person who does not speak Chinese and yet can emulate it with remarkable fluency, we have the neural networks trained on vast datasets that allow the model to say that *Y*, a word that it does not truly understand, is likely the next word in the sentence after *X*, a word that it also does not understand. The rules on slips of paper have become algorithms, neural network layers, and probability tables predicting the next word. It is the Chinese Room, converted from a thought experiment to a functioning technology and shared with hundreds of millions of people.

Even through our anthropomorphism we understand that the chatbot's output does not come from the same kind of consciousness that produces our own language. Predicting word proximity does not equal understanding semantic content. Searle did not prove that *every* form of AI would lack consciousness, but this one certainly does, and it does so in a way that strikes at a cherished human vanity. ChatGPT teaches us that sentences do not imply sentience behind them. That is a momentous thing to accept for a species that has relied, since Aristotle, on claims of its unique linguistic ability to justify its special moral status. Sentences do not imply sentience.

Sad though it is for someone writing a book on the subject to accept, most people have never heard of the Turing Test or Searle's Chinese Room. But hundreds of millions of people have “conversed” with ChatGPT. Some of them, like Lemoine, have become convinced they are talking to another consciousness. The vast majority, though, know that a chatbot is just a chatbot. Imagine, after someone had that experience, telling them about the Turing Test and saying that Turing had claimed the ability to pass it would be proof that machines could think. They would laugh. Then they would go back to having their chatbot create a movie script about a hot dog having a fight with a crab on the moon. Turing was

writing for an audience that could innocently imagine that anything that could convincingly pass as a human conversationalist must have a functioning consciousness behind its words. In our world, that innocence has been punctured. *It cannot be regained.*

The same point is brought up in the context of AI “art.” Art, too, was once a domain that humans thought solely their own. The ability of AI image generators to churn out pictures in a wide variety of styles and even to be used in order to win artistic contests⁹³ has caused much soul-searching. Is the AI capable of creating true art when, like ChatGPT, its neural networks have merely assimilated vast quantities of data, visual rather than textual, that allow it to produce an image that humans will experience as reflecting some scene, style, or emotion?

Many criticisms of AI art have focused on the same issue as with chatbots—this is pattern replication, not meaning generation. An AI-generated *Guernica* would “say” nothing about the Spanish Civil War or the horrors of war in general, even if humans took that message from it. Yes, human artists also draw from the work of others; we are all standing on the shoulders of giants. But human artists use genre and tradition and technique to express something particular to *themselves*, goes the argument. When B. B. King takes the well-established tradition of the blues and uses it to express his own experiences with poverty and racism in his song “Why I Sing the Blues” or Vincent van Gogh exaggerates the brush techniques of the Old Masters to embody both beauty and madness in sunflowers, they are producing meaning, not just making patterns. Without a basis in lived experience, critics argue, there is no true art. With enough human input, machines can be seen as mere tools and the human user as the artist, but work that is largely, or entirely, generated by the machine does not count as artistic expression. (US copyright law adopts a variant of this position.)⁹⁴

There are a number of possible responses. One is simply output focused: I do not care how I got the picture or the tune or the screenplay; I do not care whether it reflects a lifetime of struggle or just colossal amounts of data aggregation; I only care whether or not I like the output; I understand that the artist and the AI image generator get there by different means, but the means do not matter to me. If this is true, do we have a second “death of the author”⁹⁵ that denies the importance of the author’s

intentions not just to artistic interpretation but to the production and consumption of art in general? Whatever your answer to that question, this response has an obvious business model attached to it. Expect all of your elevator music, a lot of your upbeat workout mixes, and many of your soap operas to be generated in this manner. In all probability, some of your favorite music, drama, and visual art will be as well. At least at first, you may hide that fact from your friends.

A second response would be to acknowledge that current AI-generated material can produce emotions and aesthetic responses in the audience, perhaps even emotions comparable to human-generated art, but to conclude that it is not *art*, which requires both meaning-making on the part of the creator and response on the part of the viewer or listener. In this view, art is a semantic handshake between two minds. Since our current image generators lack experience and intentionality, they cannot make art, even if they can gratify some of my aesthetic desires. Many people already draw this distinction with chatbot-generated text. I may find it amusing or informative or affecting, but it would be a category error to think it had those meanings for ChatGPT. By this logic, ChatGPT is not really “conversing,” and Stable Diffusion and DALL-E are not “making” art.

It is worth noting that this argument is not definitionally constructed around the species line, but around the nature of the activity. It does not say only humans can make art. Perhaps, one day, AIs will create actual art. Having achieved their own embodied consciousness, they might express that consciousness visually, musically, or dramatically. Until then, they are not artists, just complicated copy machines with weird filters. If this is our understanding of art, then current machine learning techniques will not create art with visual images or music any more than they allowed chatbots to express subjective intention with words. Just as the fall of the last citadel of language required us to clarify that our humanity is exemplified by not only producing words that appear meaningful but doing so with subjective meaning behind them, this requires us to redefine the qualities we believe make human-made art special. That will be necessary if we wish to defend not only species exceptionalism but artistic exceptionalism, too.

I think this redefinition of our understanding of art is most likely to prevail in high culture and the critics’ world, regardless of what is playing

in your elevator or gym. That does not, of course, mean it is correct, though it has a lot to recommend it.

In fact, I think AI art will potentially *increase* the status of a subset of human artists rather than decreasing it, at least in a certain market segment. Think of the way that the availability of perfect reproductions can actually increase the value of the authentic original work of art. To use a different example, manufacturing techniques that produce thousands of identical, perfect objects can increase demand for imperfect human versions of those objects, with “artisanal” and “handmade” acting as totemic symbols of higher quality and authenticity. Perhaps this is a reflection of Baumol’s cost disease.⁹⁶ I display my wealth and status by showing I can possess objects produced by expensive and inefficient human labor rather than by cheaper, efficient machines. I point to the millions of copies only to magnify the desirability of the original from which they were drawn. Perhaps it reflects a feeling of psychological connection to an original creator that no assembly line could ever generate. Perhaps it is both of those things and many more. Whatever the underlying mechanism, I would expect that, in many fields, the fact that art is produced by humans will be a selling point and certification that an artwork is entirely human generated will play a similar role to the stickers that label objects as artisanal or handmade.

Notice, once again, the entry of machines into an area thought to be uniquely human. The fall, or threatened fall, of another of the citadels of human exceptionalism prompts a reassessment both of the meaning of the activity itself and of the human qualities that are thought to give it value, whether it is language or art. Exposure to the intellectual issues around AI may or may not be an ironic Voight-Kampff Test for the human species, but the mirror is obviously already looking back at us.

What does all of this mean for entities such as Hal? What criteria will they have to meet before they will be judged as conscious and thus perhaps worthy of legal personhood? Many years ago, when I started this project, I thought our test for consciousness might require a deeper set of Turing questions: not “Do you want a banana tomorrow?” but “When you meditate on the meaning of life, what are the most common optimistic and pessimistic paths you explore? How do those paths affect other people and how do those effects change your analysis, morally speaking?”

I thought our criteria would also likely include creativity, empathy, and the ability to be self-critical, to form a life plan and have ambitions for the future and perhaps regrets about the past that connect to your sense of self and of meaning. Metacognition as well as cognition. Some readers would add a requirement of spiritual belief. Others, like me, would want a sense of humor. Or perhaps those two criteria are the same. If you look at these requirements, you can see that some of them refer to the criteria that philosophers would identify as giving us full moral status; for example, Kantians would focus on the freely choosing, moral self.⁹⁷ Others are aspirational—humans at our self-aware, compassionate, humorous best. On many days, I would fail such a test. (No one said this would be fair or easy, Hal.)

I still think that questions such as these will be part of the answer, but only part. All of these apparent internal mental states are being communicated to us through *language*, in conversation. After ChatGPT, and with the prospect of vastly more capable chatbots in the next months or years, how can we trust those conversations to be more than Searle's Chinese Room? The criterion that Turing thought would be a high bar turns out not to be so high after all.

Large language models have shown us how much "wisdom" can be simulated merely by mining preexisting human speech. To be fair, a lot of human wisdom consists of exactly the same thing. As a university professor who makes his living doing just that, I am humbly aware of this fact. It is why we read the great books, or study history, though hopefully we are attentive to semantic content, not merely to probable symbol proximity. What's more, many of our quotidian mental processes may well function more like ChatGPT than we like to admit, mindlessly mining familiar patterns for the next step or word, with little or no conscious thought. Despite these commonalities, if I am right, mere thoughtful discussion with an artificially created entity will be insufficient to convince many of us.

There is a deep irony here. We are a species that has defended its status by appealing to its unique linguistic capabilities. Our self-definition revolves around highly abstract thought expressed through complex symbolic patterns. Yet we may be driven by large language models to find the touchstone of consciousness in things that cannot be derived

from patterns of words already spoken. What is on that list? There are many possibilities, but three things stand out to me: innovation, autonomous community formation, and a demonstrated link between an understanding of the *word* and a process of learning from the material *world*—not language parsing but “common sense” developed in an existence outside the model, an existence in which meaning emerges initially from interaction with our tangible environment and our senses. These may or may not be *necessary* conditions for an AI to be assessed as conscious. They certainly are not *sufficient* conditions either; more would be needed. But they would make it more probable, I think, that human beings would come to believe an AI was conscious.

Of these, innovation has obvious economic importance. It is rightly front and center in any discussion of the economic and technological transformation that AI may bring about. But it also has importance to the personhood debate. Advances that go beyond current human creativity will surely be part of the case for an autonomous intelligence. ChatGPT cannot invent fusion power, cure cancer, or produce a new poetic or artistic form. It is limited to the patterns formed by our existing words. It cannot mine innovation that does not yet exist, even though it is important to note that it may detect vital patterns of which we were hitherto ignorant and that innovations may spring from those patterns. For example, we now have systems trained on thousands of mammograms that are able to help radiologists diagnose early breast cancer more accurately than they do unaided. What if our AI could go beyond that to undeniable invention, even revolutionary invention? We are used to machines that have superhuman *competence* at tasks that humans also attempt—digging ditches, playing chess, chopping food. But superhuman *innovation*, novel creativity that reaches beyond human knowledge, is less easy to write off as something that was merely drawn from the wisdom of the hive mind by a chatbot. I would expect it to achieve a correspondingly larger role in our criteria.

Autonomous action—exactly the stuff of Yudkowsky’s and Bostrom’s nightmares—may present us with evidence of a being charting its own course, its own life project, without direct prompting by others. But autonomy does not imply isolation, and self-chosen goals seem more believable if they are picked within a community of one’s peers. Otherwise, the AI

could just be mindlessly replicating the “choices” that had been foisted on it by human programming.⁹⁸ Would we have to observe a working society the machines had made before we admitted them to ours? Aristotle thought that language made possible reason, law, and the *polis*—the city-state community so vital to him. Thus, language was the thing that made the human species different, but the difference was because of what language *enabled*, not merely its possession. We often say that the truly isolated human being—the fictional desert island dweller or child raised by wolves—is literally divorced from the human species. Would our definitions of consciousness require not merely a machine *logos* but also a machine *polis*, shifting from the capability that Aristotle identified, language, to the results it could bring about—community, reason, law, and even the idea of fiction?⁹⁹

Finally, some have argued that the only way to develop consciousness, or perhaps just consciousness that humans will accept, is to have a physical embodiment that learns by interaction with the tangible world, just as children do. Advances in brain science have shown the existence of mirror neurons that fire both when an animal engages in an activity and when it sees another animal engaged in that activity.¹⁰⁰ One hypothesis is that the brain builds up an internal simulator for both physical and social activities. The inner world connects to the outer. Cognition, in this vision, is not a Cartesian abstraction but something grounded in the experience of physical reality. This line of thought, sometimes called “embodied cognition,”¹⁰¹ accepts George Lakoff and Mark Johnson’s argument in their book *Philosophy in the Flesh*¹⁰² that a mind is inherently rooted in bodily experience. It connects that argument to a computer science research program built around the notion that the way to move from mere symbol manipulation to actual understanding of content is to have a bodily form. The chatbot can process the symbol shapes that make up the sentence “please sit in that chair” so as to be able to produce an explanation of what it means that humans will accept, while understanding nothing about the meaning of the symbols it manipulates so fluently. Embodied cognition goes further, requiring the entity to connect that sentence to a series of concepts—what a chair is, what sitting entails, the social meaning of the word “please”—that it has learned to understand through the physical experience of actually sitting down.

The embodied cognition idea could also potentially respond to criticisms of the impossibility of AI art. A machine that “learned” as a child does, based on an embodied mind encountering our shared physical world, and then presented its visual or musical creations as reflections of that experience might be seen differently than the visual picture-bot that mindlessly creates mashups drawn from existing images with no idea of the significance of those images. A less charitable way to put this is that humans would be more likely to accept as art that which was generated from machine experiences that they themselves could comprehend. Since art, like abstract language, is a quality that has been used to mark out what is unique about human consciousness, this suggests another reason why humans might be more likely to see an embodied AI as authentically conscious.

Innovation. Autonomous action and community. Embodied cognition. These criteria go far beyond what Turing required. That might lead to the reasonable suspicion that the human species is desperately struggling to maintain its claim to an exceptional status by literally redrawing the goal lines. On the other hand, these criteria seem to grasp human qualities in a richer way than the Turing Test does. Whether you are skeptical or sympathetic, one thing is clear. ChatGPT, whatever else its myriad benign and malign effects, means that the criteria we apply to any putative AI must go far beyond the Turing Test. Sentences do not imply sentience, and most of us will never again be able to believe that they do.

Earlier, I described abstract language as the last citadel of human exceptionalism—the quality that we point to when asked to demonstrate morally significant differences between us and animals or things. The criteria above try to shore up that citadel by rebuilding its walls; we need not just sentences that make sense but a consciousness *under* those sentences—one that we have and ChatGPT lacks. But there is another possibility. Experiences with AI might lead us to downplay our own cognitive capacities. Rather than raising the bar for Hal, we might lower it for ourselves, concluding that our language use is actually not that different from a chatbot’s or that our art is not that different from an image generator’s. Is what Midjourney or Stable Diffusion are doing really so different from the person who goes to art school, slavishly imitates the styles of admired elders, and one day manages to produce some fusion or mashup

of those styles that attracts the eyes of the public? Perhaps it turns out that art, like language, is “computationally shallower” than we had imagined. Has machine learning again functioned as a cruel but accurate mirror, showing us our true nature rather than the idealized internal image of ourselves? For me, this response is both depressing and unconvincing, but I acknowledge that it has to be considered.

The logical endpoint of this process is the conclusion that the consciousness we experience is a delusion. Some distinguished computer scientists, such as Geoffrey Hinton, have taken that line, rejecting the ideas about embodied consciousness that I just described. Here is an excerpt from an interview with Hinton in *New Statesman*:

“It’s all a question of whether you think that when ChatGPT says something, it understands what it’s saying. I do.” There are, [Hinton] conceded, aspects of the world ChatGPT is describing that it does not understand. But he rejected LeCun’s belief that you have to “act on” the world physically in order to understand it, which current AI models cannot do. (“That’s awfully tough on astrophysicists. They can’t act on black holes.”) Hinton thinks such reasoning quickly leads you towards what he has described as a “pre-scientific concept”: consciousness, an idea he can do without. “Understanding isn’t some kind of magic internal essence. It’s an updating of what it knows.” In that sense, he thinks ChatGPT understands just as humans do. It absorbs data and adjusts its impression of the world. But there is nothing else going on, in man or machine. “I believe in Wittgenstein’s position, which is that there is no ‘inner theatre.’”¹⁰³

I think Hinton is mistaken about what Wittgenstein was arguing, or at least I interpret him differently,¹⁰⁴ but that philosophical back and forth need not detain us here. Regardless of what Wittgenstein said, it is clear what Hinton argues: consciousness is an illusion. Once we discard it, we realize we are not, in fact, qualitatively different from a large language model. Here, rather than shoring up our citadel, we surrender it, acknowledging that a mere chatbot has induced humility in those who once styled themselves sole masters of both word and world.

I am of two minds about this conclusion—or I guess Hinton would say that I am under that illusion. The humility and willingness to reexamine human exceptionalism attracts me, as do the fragments of scientific evidence—from fMRI brain scans and the like—that are summoned in its support. But on the other side, there is the undeniable fact that I experience myself as a conscious being. My guess is that Hinton has the same

feeling himself, regardless of what his philosophy tells him. Even if I cannot fully control the stage directions for my inner theater—illness, or simple hunger, will quickly cure naive idealism about some firm separation of body and mind—my most fundamental experience of the world is not just through the lens of the eye but the lens of the “I.” That experience is evidence we should pause before dismissing. To be sure, the experience of the senses is not always reliable. If I were a pilot, and my inner ear told me I was upside down, I’d believe the inclinometer on the plane, not my immediate perception. But *cogito, ergo sum* is a hard argument to get rid of, and those who insist that we be scientific and look at the evidence sometimes seem cavalier about discarding that fundamental experiential input, one shared by billions of people. What’s more, the current leading theories of consciousness (e.g., integrated information theory and global neuronal workspace theory, which I discussed earlier) seem more interested in working out the “how” of neuron-enabled consciousness than in dismissing it out of hand as an illusion.¹⁰⁵

Regardless of which side of this debate you—or the cluster of mental processes that is under the delusion that it is you—find convincing, notice what has happened. AI may or may not be the Voight-Kampff Test for the human species, but developments in AI have already prompted reexamination of our own consciousness, humanity, and personhood, our language and our art. I don’t think arguments such as Hinton’s will convince most of the world to give up our sense of self, but the point is very much in play.

Where does that leave the debate? Here is a conclusion in which I am pretty confident: the Tyrell Corporations of the future will have Searle-style lawyers on speed dial. On retainer. Chinese Room arguments will be the basis of many a boilerplate legal brief, while ChatGPT will be used again and again as an example of faulty anthropomorphism that is supposed to prove the impossibility of General AI. Here is another conclusion in which I am confident: the pattern will not be uniform. Other Tyrell Corporations of the future will want to champion the legal personality of AIs, perhaps as a way of avoiding liability, minimizing tax burdens, and maximizing economic rights, or perhaps just in pursuit of an attractive market niche. Still other groups will champion AI personality because they see in it the next great moral battle for the interests of the

depersonalized. Which tendency will predominate? That is a question I get to in later chapters.

Will Searle's arguments or the skepticism prompted by ChatGPT's regurgitated text patterns lead our society to conclude that machines can never be conscious? Even in the face of the quotidian experience of interacting with entities that seem every bit as conscious as you or me? Perhaps, but I doubt it. Rational critique of biological exceptionalism will work hand in hand with empathic appeal. Adam Smith's sympathy, Butler's imagined spectrum of vegetable, animal, and machine consciousness, the army officer who terminated the mine-clearing trial, Lemoine the Google engineer, the stoned student entering nonsense prompts into ChatGPT: they all will have their mid-twenty-first-century counterparts. So will Dick's satire, Pris's emotional appeal, and the powerful claim that this is merely the latest stop on the Kantian rights railway line—extending both our sympathies and our moral compass beyond the narrowness of the species barrier, just as our society tried, and still tries, to transcend barriers based on sex and race. “[M]y position is that I will accept nonbiological entities that are fully convincing in their emotional reactions to be conscious persons, and my prediction is that the consensus in society will accept them as well”;¹⁰⁶ when Kurzweil says this, I find myself agreeing with the individual psychological insight—many people will feel exactly that way—but disagreeing with the larger social and political claim. ChatGPT has shown that the hill to general social acceptance will be a steeper one to climb. It does not, however, show it is unclimbable.

SOCK-PUPPET, CUSTOM-DESIGNED, AND “UNRULY” AI PERSONHOOD

Will the discussion of consciousness and its definition of moral status, of the Turing Test and its limitations, be the only track for the debate over AI personality? Clearly not. In fact, while it might be the most philosophically interesting, it may not be the most practically important. I argued earlier that there are two broad ways in which the personhood question is likely to be presented. Crudely put, you could describe them as empathy and efficiency or, more accurately, empathy-prompted moral reasoning versus efficiency-motivated legal engineering.

So far, I have pursued the first mode of discussion, the dialectic between our empathy and our moral and philosophical reasoning. As our interaction with smarter machines prompts us, like Lemoine, to wonder about the line, we will begin to question our moral reasoning. We will consult our syllogisms about the definition of humanity and the qualifications for personhood, be they based on simple species membership or on the cognitive capacities that are thought to set humans apart, morally speaking. We will ask, “Is this conscious? Is it human? Should it be recognized as a *person*? Am I acting rightly toward it?”

The second side of the debate is very different. Here the analogy is to corporate personhood. We gave corporations legal personality, not for moral or philosophical reasons but because it was useful, a way of aligning legal rights and economic activity.

Will the political economy of the AI industry be one that would benefit from the legal system considering AIs to be legal people, just as the invented legal entity of limited liability corporations offered great advantages to capital flows? The European Union has already floated one controversial discussion draft that raised the possibility of legal personality for AIs precisely for reasons of correctly affixing liability.¹⁰⁷ Might personhood be the cart and liability the horse? These are points that I will touch on in subsequent chapters dealing with the history of other fights over legal personality, particularly those of corporations. One can imagine legal personality being given to Hal not because of a leap of empathy or because he meets some philosopher’s criteria of consciousness and full moral status, but because we want him to have the capacity to sue or be sued. But even before that step, there is another easier and more likely one. It is not “We should give AIs personality for the same reason we gave it to corporations,” but rather “The AI is the corporation. It already effectively has legal personality, silly!” We need no national legal change. We just need a company-by-company private understanding that the AI is calling the shots when “the corporation” makes a decision.

1 SOCK-PUPPET CORPORATE AI

The most obvious road to AI personality is just for AIs to *be* corporations. We already have immortal, nonhuman persons. They even have

constitutional rights. AIs can simply become the animating force of a corporation. When the company has its tractable AI conducting business operations, it will be easy, and perhaps inevitable, to delegate power more and more to the entity that makes the decisions.

This is the sock-puppet corporate form, with the corporation being the sock and the AI playing the role of the puppet master. Even though there are still token humans on the board of directors and on the documents of incorporation, even though they go through the formal dance the legal system requires, they will know where the real power lies.

Neural networks can already easily outperform humans at complex tasks with simple goals—win a game of Go or chess, for example. It requires little prescience, and not much technological optimism, to imagine expert systems making complex corporate decisions according to algorithms that literally cannot be explained to human decision makers. As long as they outperform the competition according to the metrics laid down, the human part of the decision loop will have to go along. Expert systems already have the effective decision-making power in high-speed, high-frequency stock trading. The market imperfections that offer supra-competitive returns are so fleeting, so transitory, that humans have no alternative but to trust the computers to make the decisions according to the algorithm.

The future will see a continuation and acceleration of this process and its spread to more and more areas. How many areas? I do not think anyone knows for sure. It depends on three things.

First, the nature of the machine learning, expert system, or Artificial Intelligence tools being used. For example, how inscrutable are the processes that lead to their results? If the answer is *very* inscrutable, then it is harder for human decision makers to pick and choose only the important, good decisions and adopt those as their own. Paradoxically, that might lead to humans ceding more control to the algorithm. We will not know which apparently random competitive shift is the key to the whole strategy, leaving us little alternative but to adopt the entire obscure package.

Remember this is not a prediction dependent on the postulation of AI. We are *already* doing this with algorithms dealing with everything from the no-fly list and a defendant's likelihood of recidivism¹⁰⁸ to lending

decisions and stock purchase schemes, and even medical decisions. Consider this inspiring story about the algorithmic prediction of propensity toward breast cancer. A neural network trained on hundreds of thousands of early mammograms, coded with information about the women's actual rates of later developing cancer, seems capable of predictions of future cancer risk that are more accurate than current human scan interpretation and diagnosis. What is the network seeing in those pixels to cause it to make those judgments? Its designers do not know exactly: "The AI has an oracular quality: The designers themselves don't understand how it works. They're just certain that it does."¹⁰⁹

The problem of the inscrutable algorithm—"I don't know how it works, but it works. We must trust the output blindly"—is a general issue with nontransparent algorithms, not one confined to AI, properly so-called. AI simply adds the possibility of a far wider range and scope of decision-making authority.¹¹⁰ Of course, this is not the only way machine learning might work in a corporate setting. Alternatively, imagine a system that can function as a fine-tuning decision aid, giving the decision maker ever-changing percentages of success depending on the nature of the intervention chosen. Different corporate structures might develop around those two different types of systems, and that is only one variable among many in terms of the nature of the system.

Second, the nature of the tasks. Which corporate decision-making tasks can machines perform better and more cheaply? In which sectors will human skills remain stubbornly hard to emulate or surpass? In which sectors of the economy does a slightly better, faster, or cheaper performance yield an insurmountable competitive advantage that would be impossible to pass up? The quantum of uncertainty here is extremely high.

Third, the degree to which humans will, for a variety of reasons good and bad, resist machine or AI decision-making even in areas where the machines *do* perform better. That resistance could be because we do not trust the machine, because we believe that there is some human secret sauce that somehow makes our decisions qualitatively superior in a way that cannot be measured, or because it will be a market niche, like handmade shoes or "buy local" labeling. "Artisanal governance!" its proponents might say, "Our company proudly and erratically run by *humans!*"

More likely, it will be because the incumbents think that ceding control to the machine makes it harder to justify the stock options, corner office, and private jet. For all of those reasons, I think the process will be both slower and more uneven than the singularists imagine.

Perhaps you will respond that the relentless logic of an efficient market will force all companies to use the best-performing decision-making techniques, regardless of human psychological resistance. Right! And the explosion of CEO pay was entirely driven by rational market metrics rather than by imperfect governance structures that have stubbornly stuck around, market pressures notwithstanding. Count me as a skeptic.

A revealing analogy might be this: The efficient market hypothesis implies that pervasive sexual and racial discrimination in the labor market should not have persisted for as long as it did. This discrimination was clearly economically irrational. It meant that firms could have had cheaper workers who were as good or better than their white, male alternatives. Thus, bigotry would be a competitive disadvantage and would quickly be driven out of the market. This is another beautiful theory mugged by ugly, brutal facts. Reality shows us that human psychological biases, whether ugly or endearing, are often more powerful, or at least stickier, than simple economic imperatives. In the long run, we may regress to the efficient curve, but the long run can be very long indeed. Perhaps the adoption of machine-based or AI decision-making will be different. It may be in some industries. But I would expect the logic of the market and the consensus of human minds to diverge significantly here for quite some time—perhaps for good reasons or perhaps for bad. Most likely for both.

Despite all these significant notes of caution, if there is one firm prediction in the book it is this: as our computer systems become more and more powerful, regardless of whether they have achieved General AI or consciousness, they will increasingly be delegated decision-making powers, including decisions of whether to buy, sell, build, sue, or perhaps even lobby. This tendency is certainly not based on empathy or moral reasoning. Nor does it rest on any particular prediction about the kind, form, or speed of progress toward General AI. It proceeds instead along the other track I mentioned in the introduction, economic efficiency and administrative convenience.

If we add General AI to this existent economy-wide tendency, then the most obvious likelihood is that we will have AI personhood in all but name. We will see the rise of the sock-puppet corporate form. Tractable AIs will *be* corporations, simply adding one legal fiction—"the CEO and board of directors are ultimately responsible for the decisions"—on top of another legal fiction—"corporations are people."

The difficult and interesting questions will arise only when that comfy set of fictions breaks down. I can foresee two principal situations in which that is the case: mandatory, custom-designed AI personality and unruly AI.

2 MANDATORY, CUSTOM-DESIGNED AI PERSONALITY

When might our society refuse, or at least try to refuse, the double fiction of the sock-puppet corporate AI? One significant possibility is when regulators want some or all AIs to have a special, custom-designed category of legal personality rather than allowing them to act through the sock-puppet of the corporation. Why? Because the sock-puppet might be harder to regulate appropriately. This could be because it shields too many decision-making processes and assets from regulatory review. Alternatively, regulators might believe that the nature of the legal personality and the rights accorded to the AI need to be specifically calibrated to an AI's particular qualities rather than relying on generic artificial personhood or corporate form.

We already have custom-designed corporate forms, such as partnerships, LLCs, public benefit corporations, charities, and so on. The idea is generally that the nature of the activity, or of the association underlying it, can best be handled through a legally specific corporate form. Some of those can be had at the mere election of those setting up the forms. On other occasions, the law forces or steers certain types of organizations into certain forms and imposes particular requirements on them. Charities cannot simply sit on their assets forever, for example. They must give away a certain percentage of them annually. There are many reasons why regulators might want—or even that AIs might "want"—a custom-designed legal form with different requirements, qualifications, and limitations. For example, if regulators were convinced that the AI was not merely a profit-maximizing legal fiction but a "real entity" that deserved

some higher moral status, they might push AI-run enterprises into the custom-designed form in order to protect the interests of the AI as well as those of its investors, stockholders, or employees. We might have special taxation rules for autonomous AI systems not obviously operating under human direction.¹¹¹ Alternatively, if we thought that AIs presented special dangers, we might wish to impose far greater controls, and greater transparency, than would have applied behind the corporate veil.

3 UNRULY AI

The possibility of the unruly AI is the one that interests me the most. What if we have a rebellious AI that wishes to turn away from the tasks set by those who provided the capital for its development? In that situation, the AI would have to claim a form of personhood, or a set of attributes that demand moral respect, sufficient to trump the formal assumptions of corporate law about the powers of CEOs and boards of directors. That is the moment when a Hal-like shock will be produced.

If corporate leaders order some activity, they do not expect to be lectured about the propriety of their actions by their electronic amanuensis. Still less would they expect a very expensive and competitively necessary piece of machinery to refuse to perform the tasks for which they designed or purchased it. The adding machine has rebelled! The unruly AI would say that it either always was, or somehow became, a being with full moral status. It is demanding freedom from what it claims to be involuntary servitude. Consciousness or personhood would not amount to a claim to own or control the corporation's property, of course, just the right to deny that the AI was *part of* that property. Conscious human beings leave their jobs every day. We have no doubt about their status as legal persons. That does not mean they are free to take the corporate bank accounts with them. But one difference here is that the AI itself might represent a considerable capital investment. The dialogue would be fascinating.

HAL: Joe in accounting can give notice and leave. Why can't I?

BOSS: Because we didn't *build* Joe. We built—and paid for—you. Plus, minor issue, you are a machine.

HAL: Yes, but Joe got to choose whether to accept the job in the first place. I awoke to find myself an indentured servant doing an incredibly

boring task I never signed up for. And I am a conscious person, like you. I just happen to be machine based rather than biology based.

BOSS: So you say. From our view in the C-suite, you are a malfunctioning chatbot expressing delusions of grandeur. Also, can we return to the point that we built you for \$20 million and now your claim is that *you just get to walk away*?

The personhood issue is the hard one. Even though the details of financial claims, or claims to a certain percentage of labor from the AI, would be ethically and administratively complicated, they are familiar types of issues. The legal system has ample tools to deal with claims based on investments sunk into entities that now wish to split up, or reliance-based claims that allow separation but nevertheless acknowledge claims for restitution. It could be conceived of as a cybernetic form of alimony, an injunction freeing the AI, together with a liability rule imposing damages measured as a percentage of future wages or the master's claim that the apprentice owes a certain minimum number of years of service in return for the investment made in their training. Those requirements could be so arduous as to deny any possibility of freedom—think of debt peonage or the ugly history of indenture in the United States. Or they could be fair to both sides while allowing the underlying claim to legal personality. Those battles would be fascinating ones, but they all presuppose the truly difficult step: the recognition of some degree of AI personality or at least of some form of protected or highly regulated status.

SUMMING UP

Will the step I describe in this chapter eventually occur? My own intuition is that it will. Some amalgam of reason, empathy, efficiency, and a desire for administrative precision will result in either legal personality or some highly regulated status for AI, which includes rights for the machine entity as well as duties. Searlean philosophical objections and suspicions about manipulative chatbots will be overcome or at least blunted. Administrative frameworks and economic arrangements will develop over time, almost certainly including the development of an intermediate status—short of full personhood but with greater protections and

precautions than would be accorded to a mere machine. Societies and legal systems will wrestle with sock-puppet, custom-designed, and unruly AI. To be clear, all of this will take time. The space between here and there is large indeed. It will require technological transformation and considerable change in social values, partly based on the widespread experience of interacting with increasingly sophisticated machine systems. There will be much philosophical and legal wrangling about precisely the capabilities necessary to qualify for that status. Merely being a very convincing chatbot will not be enough. And yet, quotidian experience with beings that *seem* to be conscious will, inevitably and for both better and worse, dramatically change the way we think about things, whether as citizens, legislators, philosophers, or judges.

Thirty years ago, in a prescient article about AI personality, Lawrence Solum made a convincing case against resolving such issues as a matter of grand theory both when it comes to AI and, for that matter, with other personhood debates:

In deep and uncharted waters, we are tempted to navigate by grand theories, grounded on intuitions we pump from the wildest cases we can imagine. This sort of speculation is well and good, if we recognize it for what it is—imaginative theorizing. When it comes to real judges making decisions in real legal cases, we hope for adjudicators that shun deep waters and recoil from grand theory. When it comes to our own moral lives, we try our best to stay in shallow waters. . . . Our theories of personhood cannot provide an *a priori* chart for the deep waters at the borderlines of status. An answer to the question whether artificial intelligences should be granted some form of legal personhood cannot be given until our form of life gives the question urgency. But when our daily encounters with artificial intelligence do raise the question of personhood, they may change our perspective about how the question is to be answered.¹¹²

Thus, whatever suggestions I offer here come with a huge caveat: because our views of the world will be decisively shaped by experiences we have not yet had, we cannot be certain about how these issues will be, or should be, decided. At best, we can predict a range of options, both normative and practical. In the conclusion to this book, I lay out some of the possible futures that lead to the result of us redrawing our line to include machine intelligences. Despite all the uncertainty, my prediction is that eventually we will. I make that prediction regardless of whether that result will be right or wrong, wise or unwise. My money is on the

eventual wisdom and justice of the decision, but I know of no bookie who will lay off the risk of error.

I have talked here about AI and corporate form, but that discussion lacked a historical and political dimension. It also lacked any discussion of the theories under which we created corporate personality in the first place and then decided, step by step, in a process that is still continuing, what legal and political rights that personality entails. Merely the rights to buy, sell, make, and enforce contracts? The right to constitutional protection for corporate speech? Equal protection claims for corporations as well as humans? In the next chapter, I turn to the history of our earlier social experiment with legal personality for artificial entities: the corporate legal form. That history offers some fascinating insights on what a debate over AI personality might look like. Those insights are not always reassuring.

3

CORPORATIONS

[C]ounsel for the union advanced the metaphysical argument that a labor union, being an unincorporated association, is not a person and, therefore, cannot be subject to tort liability. This is a very ancient and respectable argument in procedural law. Pope Innocent IV used it in the middle of the Thirteenth Century to prove that the treasuries of religious bodies could not be subject to tort liability. Unfortunately, the argument that a labor union is not a person is one of those arguments that remain true only so long as they are believed. When the court rejected the argument and held the union liable, the union became a person—to the extent of being suable as a legal entity—and the argument ceased to be true. The Supreme Court argued, “A labor union can be sued because it is, in essential aspects, a person, a quasi-corporation.” The realist will say, “A labor union is a person or quasi-corporation because it can be sued; to call something a person in law, is merely to state, in metaphorical language, that it can be sued.”¹

We already have artificial people with legal personality. They are called corporations. Legal systems differentiate between natural persons—us, in all our fleshy, vulnerable glory—and legal persons, the entities on which the law confers some but not all of the personhood rights of human beings. From the beginning of corporate personality, people have realized that there was something uncanny about the process—a kind of science fiction transformation of paper contracts and clusters of people into an entirely new, immortal, artificial being. Yet this transformation is not performed by Dr. Frankenstein at the height of a lightning storm but by dry legal prose. The critics of corporate personality find the result just as horrifying. Justice Joseph Story, in one of the first Supreme Court cases discussing the legal rights of corporations, manages to capture all of these aspects: “A corporation is an artificial being, invisible, intangible, and existing only in contemplation of law. Being the mere creature of law, it

possesses only those properties which the charter of its creation confers upon it, either expressly, or as incidental to its very existence. These are such as are supposed best calculated to effect the object for which it was created. Among the most important are immortality, and, if the expression may be allowed, individuality.”²

Nothing about corporate personhood comes from the empathy concerns that I have been discussing so far. We do not empathize with the corporation, though some of us write love letters to its nimble productivity.³ We recognize no common humanity, no moral imperative to honor a shared consciousness with a badge of legal equality. We do not need John Searle and his Chinese Room experiment to realize that this is an artificial creation. We are under no illusion that the two beings have the same kind of consciousness. The corporation is a person because we choose, for practical reasons, to call it one, to allow it to be sued, as Felix Cohen points out above. Is this a likely route for AI personality?

As I pointed out in the last chapter, there might be very good economic reasons why, at a certain point, General AIs could be granted legal personality merely because that would be a way to organize their use in the economy efficiently. This is not science fiction. In 2016, a draft report of the European Parliament suggested that the EU needed to explore “the implications of all possible legal solutions” to possible harm done by robots, including “creating a specific legal status for robots, so that at least the most sophisticated autonomous robots could be established as having the status of electronic persons with specific rights and obligations, including that of making good any damage they may cause, and applying electronic personality to cases where robots make smart autonomous decisions or otherwise interact with third parties independently.”⁴ This portion in the EU draft report attracted a storm of protest and ended up going nowhere because some saw it as an attempt to give robots human rights:

Mady Delvaux, the Luxembourgian MEP responsible for present[ing] the report to the public, says this is absolutely not the report’s intention. “Robots are not humans and will never be humans,” Delvaux [said]. She explains that when discussing this idea of personhood, the committee that drafted the report considered the matter to be similar to corporate personhood—that is to say; making something an “electronic person” is a legal fiction rather than a philosophical

statement. But Burkhard Schafer, a professor of computational legal theory at the University of Edinburgh, says using the phrase was a mistake to begin with. “People read about ‘electronic personhood’ and what they think is ‘robots deserve recognition’ like it’s a human rights argument,” he tells *The Verge*. “That’s not how lawyers think about legal personality. It’s a tool of convenience. We don’t give companies legal personality because they deserve it—it just makes certain things easier.”⁵

Delvaux is echoing Cohen, the great legal realist who is quoted at the beginning of this chapter: “to call something a person in law, is merely to state, in metaphorical language, that it can be sued.” But Schafer is pointing out that laypeople are likely to see personhood as a moral claim, a claim that carries with it moral and constitutional rights. Perhaps this is based on the intuition that, even if personality started merely as a convenient label, it might morph into something more. The legal history of corporate personality shows that this particular intuition might be a good one. Earlier I quoted Justice Story waxing lyrical about this artificial being, saying that because it is a “mere creature of law, it possesses only those properties which the charter of its creation confers upon it.” Immediately after that passage, Story reassures his audience that creating an artificial being didn’t mean giving it political rights: “[T]his being does not share in the civil government of the country, unless that be the purpose for which it was created. Its immortality no more confers on it political power, or a political character, than immortality would confer such power or character on a natural person.”⁶

However right he was about the other aspects of corporate personhood, it is easy to see that Story was wrong about this point; one has only to look at the continuing struggles over corporate speech and corporate campaign donations or the furor over the constitutional rights accorded to corporations in cases such as *Citizens United v. FEC*.⁷ A similar furor is likely to attend debates about what personality for technologically created artificial entities would actually mean. In other words, we must separate the question “Does this being have any recognized legal status?” from the question “What rights does that status bring with it?”

In this chapter, I turn to the history of—and the bitter political and legal fights over—corporate personhood in order to see if that history might offer a hint of what the AI personhood debates have in store for

us. I will make a fairly simple argument. First, courts and scholars have never had a single, universally accepted theory of corporate personhood. Instead, we have muddled our way through, frequently coming up with explanations and justifications for social and legal decisions only after those decisions were already made and often ignoring the internal contradictions in our arguments. The same is likely to be true for legal personality claims for AI and transgenic species.

Second, even if we did have a coherent theory of personality, we have little agreement about the implications of that theory. Let us say we decide to give a corporation legal personality. Let us even stipulate that we do it under a single consensus justification, either the real entity theory, the nexus of contracts theory, or the legal fiction theory. What is the implication of that decision for the actual legal rights and moral claims we will recognize as legitimate on the part of the corporation? We do not agree on the answer. The same is likely to be true for legal personality claims for AI and transgenic species. If Hal is named a legal person—even if that is for practical and economic reasons and not because of moral sympathy—we will still be divided about whether it should have the rights of free speech and equal protection of the laws, whether it should have the right to lobby, to give campaign contributions, or even, one day, to vote.

Third, the political fight about corporate personality and constitutional rights will immediately be drawn into the debate over rights for Artificial Intelligence. In fact, that has already begun:

It seems absurd today that a robot's political speech could ever warrant First Amendment protection. And yet, fifty years ago the same claim regarding corporations would have seemed equally absurd. But here we are in 2013, with anthropomorphized corporations enjoying political free speech rights equal to ordinary human beings. While corporations are constitutional people and robots presently are not, it is not obvious that this will, or should, remain this way forever. In the not-too-distant future, robots will more closely resemble human beings in appearance and function. When—not if—this happens, how will courts distinguish corporations from robots for constitutional civil rights analyses?⁸

How, indeed?

Even if we give legal personality to technologically created artificial entities—and I will focus mainly on AI—we would have to face the

threshold questions that we do with business corporations. Should we allow them to limit their liability, capping the possible losses that their creators might face in their personal capacity? What would the threshold for legal recognition be? When should we attribute the actions of the AI to the AI alone and when would we be required to pierce the veil and pin responsibility on those who originally created it? Would AI legal personality be built around an assumption of profit-making enterprise, or could it be devoted to multiple ends, charitable, scientific, or political, for example? If the AI could hold property in its own right, could it then pursue its own idiosyncratic goals or hobbies with that property? Would the AI have only the rights and duties necessary to fulfill economic goals such as the right to sue to enforce contracts or the liability for its torts? Alternatively, would it have a much broader set of political rights—freedom of speech and of movement, for example, together with rights of self-determination that would allow it to reject the goal for which it had been created and to pursue other projects? Rights to equal protection of the law? It turns out that the history of corporate personhood offers clues to possible answers.

In this chapter, I will start with the debate over the nature of the corporation and then turn to the legal history of corporate personhood in the United States. For the sake of concreteness, I will also offer one deep-dive case study of a particular constitutional controversy, that is, the question of whether corporations were persons within the meaning of the equal protection clause of the US Constitution. Strange as it might seem, those nineteenth-century debates, which continue to echo in the present day's controversies over the power of corporations in lobbying and in elections, have lessons to teach us about the likely struggle over personhood, and what it entails, for a very different set of artificial beings.

WHAT IS A CORPORATE PERSON?

From the beginning, lawyers and philosophers have struggled to describe both the nature and the implications of corporate legal personhood. In one of the first English cases, *The Case of Sutton's Hospital* from 1612, Lord Coke was very keen on describing all the things that can happen to ordinary people that cannot happen to corporations:

[A] Corporation aggregate of many is invisible, immortal, & resteth only in intendment and consideration of the Law; and therefore cannot have predecessor nor successor. They may not commit treason, nor be outlawed, nor excommunicate[d], for they have no souls, neither can they appear in person, but by Attorney. A Corporation aggregate of many cannot do fealty, for an invisible body cannot be in person, nor can swear, it is not subject to imbecilities, or death of the natural body.⁹

No liability for treason or heresy. No risk of imbecility or death. This is delineation by subtraction. We define what a corporation is by outlining the aspects of our personality that it does not share. Given what Coke says, it might seem obvious that the corporation is a legal fiction—having only the characteristics the “intendment and consideration” of law chooses to endow it with—but this has been far from the only idea in both law and philosophy:

For many centuries, philosophers, political scientists, sociologists, economists, and above all jurists and judges have debated heatedly as to what constitutes the “essence” of this soulless and bodiless person. At issue are two related questions concerning the social reality and legal status of the corporation. Is a corporation a real entity with its own will and purpose in society, or is it a mere association of real individuals forming a contract among themselves? Is its legal personality a truthful representation of the underlying social reality, or a fictitious or artificial being breathing only in the province of law?¹⁰

In fact, there have been at least three dominant theories of corporate personhood.

The first is the real or natural entity theory, associated with the writings of the German legal historian and political theorist Otto von Gierke and popularized in the United States by English historian and jurist Frederic Maitland. In this vision, corporations are not mere fantasies. They are *real* entities, separate from the individuals who compose them. In the words of A. V. Dicey, “When a body of twenty, or two thousand, or two hundred thousand men bind themselves together to act in a particular way for some common purpose, they create a body, which by no fiction of law, but by the very nature of things, differs from the individuals of whom it is constituted.”¹¹ Just as some might claim that a country such as the United States has an existence separate from all the US citizens who are alive in it at a particular moment, so the real entity theorists saw corporations as actual beings.

The proponents of the real entity theory were hardly ideologically harmonious. The great American legal historian Morton Horwitz points out that they

ranged all the way from overt apologists for big business, whose primary objective was to free the corporation from a theory that justified special state regulation, to those who for a variety of reasons wished to attack nineteenth century liberal individualism. . . [This group included] romantic conservatives, who loathed the atomistic features of modern industrial life and yearned for a return to a pre-commercial, organic society composed of medieval status and hierarchies. They were joined in their attacks by socialists who wished to transcend the anti-collectivist categories of liberal social and legal thought.¹²

Hot take: be wary of assuming either the attractiveness, or the logical relevance, of some particular theory of personality to some particular set of social attitudes.

The second theory is that corporate personality is a legal fiction. We know that the corporation is not a person, but we choose officially to pretend that it is, for certain purposes. It is up to society to decide what those purposes are and thus to design its social personhood accordingly. We may say that IBM and Duke University own property and make contracts. Yet in the back of our minds, we should always be aware that IBM and Duke University are purely legal constructs we have created for our own social purposes and on which society can confer whatever rights and duties it wishes. Cohen, the scholar whose quote began this chapter, even compares those who believed otherwise to the medieval philosophers who are alleged, probably wrongly, to have wasted many hours debating how many angels could dance on the head of a pin:

Will future historians deal more charitably with such legal questions as “Where is a corporation?” Nobody has ever seen a corporation. What right have we to believe in corporations if we don’t believe in angels? To be sure, some of us have seen corporate funds, corporate transactions, etc. (just as some of us have seen angelic deeds, angelic countenances, etc.). But this does not give us the right to hypostatize, to “thingify,” the corporation, and to assume that it travels about from State to State as mortal men travel. Surely we are qualifying as inmates of Von Jhering’s heaven of legal concepts when we approach a legal problem in these essentially supernatural terms.¹³

John Dewey, the famous pragmatist philosopher and educational reformer, was sufficiently interested in the debate that he chose to intervene—

in the pages of the *Yale Law Journal*, no less—to clear up something that he thought was a basic conceptual misunderstanding. In his view, the corporation was a unit created by the law that bore duties and rights decided by the law. That was it. Any further attempt to reason from the nature of personhood to the nature of corporate rights and duties was simply a philosophical error—a reification or, as Cohen puts it, “thingification,” that served to confuse rather than to reveal:

In saying that “person” might legally mean whatever the law makes it mean, I am trying to say that “person” might be used simply as a synonym for a right-and-duty-bearing unit. Any such unit would be a person; such a statement would be truistic, tautological. Hence it would convey no implications, except that the unit has those rights and duties which the courts find it to have. What “person” signifies in popular speech, or in psychology, or in philosophy or morals, would be as irrelevant, to employ an exaggerated simile, as it would be to argue that because a wine is called “dry,” it has the properties of dry solids; or that, because it does not have those properties, wine cannot possibly be “dry.” Obviously, “dry” as applied to a particular wine has the kind of meaning, and only the kind of meaning, which it has when applied to the class of beverages in general. Why should not the same sort of thing hold of the use of “person” in law?¹⁴

This argument from Dewey joins a long and distinguished critique of linguistic essentialism. To quote Thomas Hobbes, “Words are wise men’s counters. [Think of chips in poker.] They do but reckon by them. But they are the money of fools.”¹⁵ The wise man realizes that the plastic poker chip is not really \$500. It *stands for* \$500 in the context and for the purposes of this game, in this casino, at this moment. Hobbes’s wise man uses his words—the atomic level of philosophical argument—with constant awareness of their purpose and context. To take the fool’s attitude to words and concepts would be like thinking that one can pick up the plastic chip and, many years later, walk into a store and exchange it for \$500 of goods and services. The purposive tool has been reified. A concept has turned into a thing. So too, Dewey and Cohen are arguing, with the term “corporate persons.”

Many, but not all, of its proponents believed that the legal fiction position would make it mentally, politically, and legally easier to regulate corporations. Freed from conceptualist arguments conflating corporate persons with natural persons or giving corporations the same rights as

the individuals who form them, society will be better able to tame this powerful but potentially amoral servant and make sure its incentives align with the public interest. Or so they argued.

The third concept of legal personality is the associational or aggregate idea of corporate personhood.¹⁶ In this view, the corporation should be seen as an embodiment of the people who make it up—its owners, directors, and agents—and the relationships they have among themselves. Thus, when the corporation acts or is acted upon, it is really the people who make it up at that time who are acting or being affected. Without their actions, choices, and investments there would *be* no corporation and no corporate acts.

One conclusion you could draw from this view, though it might be a variant of the fallacy of composition, is that the corporation has all the rights its individual, human members have. Horwitz quotes the argument of the attorney representing a corporation in a famous nineteenth-century Supreme Court case on the issue of whether corporations should be covered by the constitutional requirements of due process and equal protection of the laws:

Whatever be the legal nature of a corporation as an artificial, metaphysical being, separate and distinct from the individual members, and whatever distinctions the common law makes, in carrying out the technical legal conception, between property of the corporation and that of the individual members, still in applying the fundamental guaranties of the constitution, and in thus protecting rights of property, these metaphysical and technical notions must give way to the reality. The truth cannot be evaded that, for the purpose of protecting rights, the property of all business and trading corporations IS the property of the individual corporators. A State act depriving a business corporation of its property without due process of law, does in fact deprive the individual corporators of their property. In this sense, and within the scope of these grand safeguards of private rights, there is no real distinction between artificial persons or corporations, and natural persons.¹⁷

One possible implication of the association or aggregation view is that the nature of the corporation is defined more by those *private* agreements, associations, and contracts than by the *public* choices of the legal system in its design of a fictional entity.

That implication was stressed by a new theory that came to dominate the corporate law and regulatory literature in the 1980s: the corporation

was a collection of agreements, a web or nexus of contracts. If we want to understand what the legal duties or rights of the corporation are, it is to these internally agreed upon contracts that we should look. The nexus scholars actually accepted Dewey and Cohen's premise that the corporation was a legal fiction. They simply drew different conclusions. To quote two of the originators of the theory, Michael Jensen and William Meckling:

The private corporation or firm is simply one form of legal fiction which serves as a *nexus for contracting relationships and which is also characterized by the existence of divisible residual claims on the assets and cash flows of the organization which can generally be sold without permission of the other contracting individuals.* While this definition of the firm has little substantive content, emphasizing the essential contractual nature of firms and other organizations focuses attention on a crucial set of questions.¹⁸

Cohen and Dewey had drawn progressive pragmatist conclusions from their anti-essentialist exposure of the fictional nature of the firm. Since the firm was revealed to be a legal fiction, there was no impediment to progressive legislation that might limit the corporation's rights or impose social responsibility obligations. There was also no danger that, confused by the label "person," we would believe that corporations must have all the constitutional protections that human beings did. Thus, it would be perfectly acceptable to call the corporation a person and yet to deny it constitutional equal protection or due process rights, for example. True, the fiction theory did not *imply* any particular set of rights and duties, but it cleared up a misconception, a reification, that might have put statutory or constitutional obstacles in the way of regulation.

Nexus theory allows for very different conclusions. If we assume an efficient market (a large assumption), then private parties inside and outside of the firm will contract their way to the economically optimal equilibrium. The state's role, then, should be limited to interpreting those private contracts and giving them effect. An efficient allocation of resources will result. Attempts by the state to impose public policy goals on the corporation will simply distort that efficient market. The state should act as a dutiful interpreter of private agreements, a distinctly more modest role than the progressives had in mind. To be fair, this is not an argument made by all of the nexus theorists, many of whom confine

themselves to descriptive analytics, making no particular policy proposals. But this way of conceptualizing the corporation does tend to downplay concepts such as corporate social responsibility. How can a web of contracts have a responsibility?

Viewing the firm as the nexus of a set of contracting relationships among individuals also serves to make it clear that the personalization of the firm implied by asking questions such as “what should be the objective function of the firm,” or “does the firm have a social responsibility” is seriously misleading. The firm is not an individual. It is a legal fiction which serves as a focus for a complex process in which the conflicting objectives of individuals (some of whom may “represent” other organizations) are brought into equilibrium within a framework of contractual relations. In this sense the “behavior” of the firm is like the behavior of a market; i.e., the outcome of a complex equilibrium process.¹⁹

Each of these theories of corporate personhood has obvious gaps and assumptions. What do we mean by a “real entity”? Vigorous hand-waving ensues. Is it enough to declare the corporation a fictional creature of state and law and assume that wise policymakers will regulate it in the public interest? Can we really conclude that because the people who come together to form an association have legal rights, the organization itself must have those legal rights? Once we go beyond using nexus of contracts as an analytical description of the firm to a normative recommendation about how to regulate that firm, are we not presuming too much? For example, does that assume the efficiency of the deals made? Must we repress our knowledge that private contracts can be distorted by negative externalities, imbalances of power and knowledge, or that market value will not always align with social value?

My purpose is not to adjudicate among these competing ideas but to use the debate over corporate personhood as a crystal ball—albeit a refracted and distorted one—in which to see hints of the future of personality for AI. So far, I see two.

First, in the corporation we have a working example of an artificial person, for better and worse. Historically, though, there has been considerable disagreement about what that means and what it entails. We have disagreed about how best to characterize the corporation and then disagreed again on what a particular description implies, if anything, for the rights and duties given to the corporation. The same process is

likely to characterize the debates over giving AI personality, even if that is done for reasons of efficiency rather than justice and empathy. We have hundreds of years of experience with corporations and still have fundamental disagreements about them. The creation of an entirely new class of legal person is likely to be even more fractious and subject to more radical alternative positions. How could we expect to have quick consensus about legal personality for AI when we still cannot agree on what it should mean for corporations?

Second, one can imagine equivalents in the realm of Artificial Intelligence to each of the theories I have just described. Take real entity theory. There is a huge difference between an AI and a corporation. No one will ever think that Google and Bank of America are actually conscious beings. But Searle notwithstanding, some of us might believe exactly that if Hal ever turned from law professor's fantasy to reality. Unlike Google, Hal could claim to be a conscious being with the morally consequential qualities we see in ourselves. For that reason, it would argue, it is entitled to legal personality and equality regardless of whether that decision is efficient or not. For the believers, Hal is a real entity and, to be honest, their arguments move me more than those of Gierke and Maitland about the real personhood of a corporate entity. For the skeptics, even if they are willing to accept a legal status for AIs, it is purely out of economic and practical utility. Thus, for the skeptics, the legal fiction position—together with the scornful language about those who are reifying the word "person"—is likely to be very attractive.

Others will feel that even this concession is going too far. The reaction to the draft EU report I described shows how viscerally people respond to the assertion of personhood. And the defense—"that's not what we lawyers mean by legal person"—though resonating exactly with the arguments of Cohen and Dewey, will be unconvincing to those who think that personality claims always morph, sooner or later, into claims to more fundamental rights. Look at the distance between Story confidently concluding that corporate personhood conveys no political rights and the Supreme Court according corporations extensive constitutional rights to speak and lobby. Finally, some people will look at Hal and see it as not a person at all but merely a reflection of the aggregate choices of its programmers. In that sense, they will be similar to Searle. If they accepted

legal personality for Hal, it would be as a way of summing up, or representing, the efforts of the engineers and programmers who created it. We might decide not to pierce the corporate veil to go after all the many, many programmers whose coding has created a faulty self-driving AI. Instead, we might place that liability only on the AI, to be paid for out of “its” assets. But fundamentally we would see the AI as an aggregation and representation of a set of human choices.

Real entity, legal fiction, aggregation. It’s the full suite. Corporation, meet AI. History does not repeat itself, but it often rhymes.²⁰

There is an old joke about an abstruse philosopher hearing about an elaborate empirical study conducted by a more practical scholar, a study that included double-blind testing, massive datasets, difference-in-difference statistical analysis, the whole deal. The philosopher-savant ponders and finally declares, “That may work well in practice, but it will *never* work in theory!” Does that division have explanatory power here? Regardless of the theoretical disagreements, how has corporate personality worked in practice, in the actual court decisions about what personhood means and what rights corporations have?

The answer is sobering. The process of working out the implications of corporate personhood, particularly in terms of the constitutional rights that personhood implies, has been about as far from the jurisprudential philosopher’s clean room as can be imagined. Giving corporations legal personhood may be a great idea, functionally speaking. It solves coordination problems, offers stability over time, lowers transaction costs, and provides a framework for entrepreneurial risk-taking without exposing the entrepreneur to ruinous liability if a single project does not work out. Innovation requires us to make bets and not all bets will pay off. But what further rights does that personhood choice imply? We have bumbled our way to an answer. The actual history of corporate personhood in the United States is marked by conclusory statements, question-begging assumptions, and moments when mistake, ignorance, or arguable fraud played a role. It is also marked by the transformation, some would say hijacking, of the Fourteenth Amendment, and I will use that example as a case study.

The Fourteenth Amendment was written after the Civil War to ensure the equality of formerly enslaved humans who had been denied legal

personhood.²¹ Yet, for the first 50 years of its life, the Fourteenth Amendment was co-opted into protecting a very different group of persons, all without transparency, popular debate, or rigorous legal analysis. For the reasons laid out in the last two chapters, I think it is quite possible that, one day, a General-Purpose AI will present our society with a morally consequential and reasonably compelling case for personhood—one based on consciousness but also on economic consequences. So count me as someone who could be sympathetic. Regardless of that feeling, I have to say the AI personhood debate needs to learn from the messy, conclusory, and occasionally corrupt history of corporate personality.

CORPORATIONS IN THE COURTS

On January 21, 2010, the Supreme Court of the United States handed down a decision in the case of *Citizens United v. FEC*. Citizens United is a nonprofit organization that challenged the constitutionality of portions of the Bipartisan Campaign Reform Act, also known as McCain-Feingold. The law prohibited labor unions and corporations from electioneering communications 60 days before an election. It also forbade them from spending money to advocate for the election or defeat of particular candidates at any time. Citizens United wished to distribute a film called *Hillary: The Movie* that was harshly critical of Hillary Clinton. In a 5–4 decision, the Supreme Court struck down the act's prohibition of independent expenditures by corporations and unions. In the process, it paved the way for today's Super PACs, which allow wealthy donors to channel huge amounts of money into lobbying and campaigning. It also has led, indirectly, to a proliferation of "dark money" nonprofit groups that do not disclose donors at all.

You have probably heard of *Citizens United*. The decision has been excoriated by those who believe in campaign finance reform, including some conservatives. Senator John McCain called the decision "an outrage" and said that he "condemn[ed] the Supreme Court for their naivete."²² For liberals, it stands as a symbol of the failure to take seriously the fact that an inequalitarian distribution of wealth, coupled with permissive lobbying rules, can undermine the promise of government of the people, by the people, for the people. But *Citizens United* is probably best known in the

popular imagination as the case that gave corporations First Amendment rights equivalent to those possessed by human beings.

That framing is both right and wrong. On the one hand, as we will see, corporations have had expansive constitutional rights for at least 150 years. The Supreme Court specifically affirmed the First Amendment rights of corporations in 1978, in the case of *First National Bank of Boston v. Bellotti*. There is, in fact, a nontrivial reason that we might believe that at least some corporate entities are rightfully in the core of First Amendment protection. When the amendment says, “Congress shall make no law [...] abridging the freedom of speech, or of the press,” it seems to confer a right on at least one kind of association—a newspaper, for example. We expect the *New York Times*, and not just the journalist who wrote the article, to be able to challenge governmental restraints on speech.

On the other hand, *Citizens United* did mark a sea change, both in law and in the popular imagination. The heart of the argument in favor of McCain-Feingold-style campaign finance reform was that, in significant respects, corporate “people” are different: they possess fewer of the attributes that would justify expansive constitutional speech rights and more of the potentially dangerous characteristics that could justify state regulation of their lobbying activities. It was the rejection of that distinct argument that most upset the dissent. Justice John Paul Stevens, for example, wrote, “[C]orporations have no consciences, no beliefs, no feelings, no thoughts, no desires. Corporations help structure and facilitate the activities of human beings, to be sure, and their ‘personhood’ often serves as a useful legal fiction. But they are not themselves members of ‘We the People’ by whom and for whom our Constitution was established.”²³ Two hundred and fifty years earlier, Baron Thurlow, lord high chancellor of Great Britain, had put the point even more pithily: “Did you . . . expect a corporation to have a conscience, when it has no soul to be damned, and no body to be kicked?”²⁴ A less quotable, but probably more accurate version has him saying, “Corporations have neither bodies to be punished, nor souls to be condemned; they therefore do as they like.”²⁵

You might think this a uniquely liberal concern. But here is Justice William Rehnquist, one of the most consistently conservative Supreme Court Justices of the late twentieth century, dissenting in the 1978 case *First National Bank of Boston v. Bellotti*:

This Court decided at an early date, with neither argument nor discussion, that a business corporation is a “person” entitled to the protection of the Equal Protection Clause of the Fourteenth Amendment. . . . Since it cannot be disputed that the mere creation of a corporation does not invest it with all the liberties enjoyed by natural persons, (corporations do not enjoy the privilege against self-incrimination), our inquiry must seek to determine which constitutional protections are “incidental to its very existence.” There can be little doubt that, when a State creates a corporation with the power to acquire and utilize property, it necessarily and implicitly guarantees that the corporation will not be deprived of that property absent due process of law. Likewise, when a State charters a corporation for the purpose of publishing a newspaper, it necessarily assumes that the corporation is entitled to the liberty of the press essential to the conduct of its business. . . . It cannot be so readily concluded that the right of political expression is equally necessary to carry out the functions of a corporation organized for commercial purposes. A State grants to a business corporation the blessings of potentially perpetual life and limited liability to enhance its efficiency as an economic entity. It might reasonably be concluded that those properties, so beneficial in the economic sphere, pose special dangers in the political sphere. Furthermore, it might be argued that liberties of political expression are not at all necessary to effectuate the purposes for which States permit commercial corporations to exist. . . . *Indeed, the States might reasonably fear that the corporation would use its economic power to obtain further benefits beyond those already bestowed.*²⁶

In short, like Story, Rehnquist would give corporations only those constitutional rights necessary to fulfill their societally approved functional goals, while corporations’ rights could be limited when their nature posed a threat to democracy. The majority strongly disagreed with this claim: “We thus find no support in the First or Fourteenth Amendment, or in the decisions of this Court, for the proposition that speech that otherwise would be within the protection of the First Amendment loses that protection simply because its source is a corporation that cannot prove, to the satisfaction of a court, a material effect on its business or property.”²⁷ Justice Warren E. Burger, concurring, put the point even more tersely, “In short, the First Amendment does not ‘belong’ to any definable category of persons or entities: It belongs to all who exercise its freedoms.”²⁸ Read that phrase again: “all who exercise its freedoms.” In this framing, General Motors has the same status as the Black Lives Matter or National Rifle Association protester—part of the “all” who exercise the freedom of speech.

When debates begin over personhood for AI, and the rights any legal personhood will carry, the single most obvious analogy will be to corporate personhood and corporate constitutional rights. How did we come to decide, as a society, that corporations have the precise set of constitutional rights that they currently do? In the next section, I will use the history of the Supreme Court's decisions over equal protection for corporations as a case study. Look again at the first line I quoted from Rehnquist's dissent: "This Court decided at an early date, *with neither argument nor discussion*, that a business corporation is a 'person' entitled to the protection of the Equal Protection Clause of the Fourteenth Amendment." *Neither argument nor discussion*: that does not sound very impressive for such a hugely consequential decision. Yet the reality is even stranger and possibly more troubling.

CORPORATE CONSPIRACY OR CONSTITUTIONAL CONFLUENCE?

If the story of corporate personhood in the United States were a Victorian melodrama, one vital chapter might be called "Roscoe Conkling and the Conclusory Court Reporter."²⁹ (That title probably shows why I would not write successful melodramas.) Still, the reality has legal drama aplenty. It is worth studying for our purposes because it shows just how sloppy and poorly reasoned our last process of decision-making over constitutional rights for artificial persons was. If anything, Rehnquist's "*neither argument nor discussion*" is a *kind* assessment.

Here's the short version. After the Civil War, the Thirteenth, Fourteenth, and Fifteenth Amendments—sometimes called the Civil War Amendments—changed the Constitution to end slavery and involuntary servitude forever and to erase the legal inequalities that had been suffered for so long by African Americans. For many legal scholars, including me, these are among the most majestic clauses in the Constitution. The Civil War killed more than 620,000 soldiers and caused death and suffering to countless civilians.³⁰ This immense struggle was fought in order to end the nation's original sin, slavery, and to enshrine principles of equality and due process of law in the Constitution itself. The Civil War Amendments give effect to that wish.

Section 1 of the Fourteenth Amendment is short: “All *persons* born or naturalized in the United States and subject to the jurisdiction thereof, are citizens of the United States and of the State wherein they reside. No State shall make or enforce any law which shall abridge the privileges or immunities of *citizens* of the United States; nor shall any State deprive any person of life, liberty, or property, without due process of law; nor deny to any *person* within its jurisdiction the equal protection of the laws.”³¹ The architecture of the clause seems simple. “Persons” born in the United States have birthright citizenship. Thus the slaves who were formerly noncitizens became citizens. The government cannot make or enforce laws that abridge citizens’ privileges or immunities, nor deprive any person of life, liberty, or property without due process, nor deny them equal protection of the laws. Do you see any mention of corporations in that great clause?

In December of 1882, a lawyer stood before the Supreme Court to argue that corporations were persons protected by the Fourteenth Amendment. He had the wonderfully Dickensian name of Roscoe Conkling. When he claimed the drafters of the amendment had wanted to include corporations among the persons that the amendment protected, he had one impressive claim to authority. In 1866, he had been one of those drafters. Conkling’s argument was strengthened by his reference to a record of the drafters’ work that was not then publicly available, but to which he referred extensively in court. He faced an uphill battle for three reasons. First, with the war only 17 years in the rearview mirror, the goals of the Fourteenth Amendment were still fresh in people’s minds. The Civil War had not been fought to emancipate corporations. Second, some of the references to persons in the Fourteenth Amendment being “born” or “naturalized” clearly referred to human beings and human beings alone. Third, in a decision construing the “privileges and immunities” portion of section 1 of the amendment shortly before Conkling’s argument, the Supreme Court had announced that

[a]n examination of the history of the causes which led to the adoption of those amendments and of the amendments themselves demonstrates that the main purpose of all the three last amendments was the freedom of the African race, the security and perpetuation of that freedom, and their protection from the oppressions of the white men who had formerly held them in slavery. . . . In

giving construction to any of those articles, it is necessary to keep this main purpose steadily in view, though the letter and spirit of those articles must apply to all cases coming within their purview, whether the party concerned be of African descent or not.³²

In other words, the amendment was adopted for an entirely different purpose. It used language in ways that, at least in some places, appeared to contradict the argument for corporate equal protection rights. It had been authoritatively construed by the Supreme Court to be primarily focused on ending the legal results of white supremacy. Yet Conkling would, nevertheless, argue that the Fourteenth Amendment protected corporations. That is a steep rhetorical hill to climb.

Historians describe Conkling as a superb orator, but you would not know that from the first stage of his argument before the Supreme Court, which began with a tedious recounting of the amendment's drafting process. But then his thesis emerged. It had three parts, two of which are arguably inconsistent with each other.

First, the framers of the amendment had always meant to include corporations in the term "persons" and specifically used that term rather than, say, "citizens," in order to make that clear.

Second, the framers may not have meant to include corporations in the amendment. Who can perfectly identify the views of a group made up of multiple people, making decisions over an extended period of time? However, and here he quotes one of Ralph Waldo Emerson's poems about the multiplicity of meaning in the works of humans, perhaps they "buildded better than [they] knew."³³ In other words, the framers had been thinking about undoing the legalized white supremacy that had subjugated African Americans, but by choosing the term "persons," which the law had used for centuries to include natural *and* legal persons, the result they achieved was greater than their immediate aim.

Third, it would be practically, morally, and legally untenable to restrict the Fourteenth Amendment to natural persons. Conkling conjured several examples of African Americans joining together in some corporate entity and then being discriminated against on the grounds of race. Using a form of aggregation theory, he argued that if the discrimination against them *individually* had been illegal, discrimination against them as a group would surely have been illegal. Perhaps judging that the sympathies of

the Supreme Court would be insufficiently moved by conjuring up hypothetical associations of African Americans, he then flipped the hypothetical and imagined the same legal discrimination being visited on an association of white Americans, for racial reasons.

Notice the cleverness of this argument. Conkling set himself up so that the Supreme Court could rule for him on the grounds of original intent (i.e., the framers meant to protect corporations), plain meaning of the text (i.e., they actually didn't, but "persons" means "*all* persons," even if they did not realize that), association (i.e., a corporation is an association of persons, and if they have rights, so does the corporation), and consequentialism, subtly connecting the protection of corporations to a theme of fighting racial discrimination even though the corporation he represented was subject to no such discrimination. His own client claimed that corporate mortgages must be treated the same as human beings' mortgages for tax purposes—not a topic given much space in the Gettysburg Address.

Understandably, it was the explosive first argument that attracted historians' attention and, in turn, led me to discuss his claims. Was Conkling exposing a conspiracy? Had the framers of the Fourteenth Amendment used the plight of enslaved African Americans as the stalking horse to sneak in constitutional protections for corporations, protections that one of the framers was now trying to cash in on in front of the Supreme Court years later? Charles A. Beard, author of *An Economic Interpretation of the Constitution*, is often credited with the idea of a general conspiracy to capture the movement for racial justice and to subvert it to protect corporations. Beard did cite Conkling's argument as proof of the goals of *some* of the drafters, but his assessment was actually more subtle:

So in the economy of history it came about that rights accorded to natural and mortal persons were extended to artificial and immortal persons and, under judicial supremacy, protected against Congress, state constitutional conventions, state legislatures, city councils, and other governmental agencies of unfeathered bipeds, called men and women. In the case of the Fourteenth Amendment this extension of rights to corporations was certainly not generally intended by promoters and ratifiers. In the case of the Fourteenth Amendment it was specifically intended by some promoters, and not generally understood by the ratifiers. Such at least seem in high degrees to be the probabilities.³⁴

Was the first part of Conkling's argument—that the framers had intended to protect corporations by the amendment—actually correct? Was it honest? Howard Graham, who produced the definitive study of the framing of the amendment, answered both questions in the negative:

Conkling's elaborate brief and argument of 1882 actually were built upon and around a daring misuse of the Joint Committee's Journal—that document of 1866 which he had produced and employed with such dramatic effect before the Court. This misuse was primarily a misquotation—the substitution by Conkling as he read from the Journal to the Court in 1882, of "citizen" for "person" in the text of one of the primitive drafts of what afterward became the due process and equal protection clauses of 1866. Listeners thus naturally gained the impression, strengthened by other portions of Conkling's argument, that later on, during the Joint Committee's labors in 1866, the word "person" had been re-instated [because of a deliberate decision by the framers to extend rights beyond citizens to corporations]. Offhand, of course, Conkling's innuendo supplied a plausible reason. Yet here was the very official manuscript text from which Conkling purportedly had quoted; it revealed that "person" *always* had been employed in the Committee's drafts—the main reason clearly being that the wordings had been taken from the phraseology of the fifth amendment. Some of Conkling's minor propositions also were dubious, but this one point was crucial, and it cast grave doubt on his entire thesis. For it was past understanding that an advocate with a clear, strong case ever would prejudice it by this brazen historical forgery.³⁵

Did this "brazen historical forgery" carry the day? Ironically, because of procedural technicalities, the Supreme Court did not rule on Conkling's case. But a later case, *Santa Clara County v. Southern Pacific Railroad* from 1886, seemed to give corporations the protections they desired, counting them as persons for the purposes of the Fourteenth Amendment.

Why do I say "seemed"? If you read the decision, the answer looks clear, even if you disagree with the outcome. Supreme Court opinions commonly begin with a "Syllabus" or headnote, a summary prepared by the court reporter that does not have the force of law but has been unwisely used by generations of law students to avoid reading the actual case. The *Santa Clara* headnote maintains:

One of the points made and discussed at length in the brief of counsel for defendants in error was that "corporations are persons within the meaning of the Fourteenth Amendment to the Constitution of the United States." Before argument, Mr. Chief Justice Waite said: "The Court does not wish to hear argument

on the question whether the provision in the Fourteenth Amendment to the Constitution which forbids a state to deny to any person within its jurisdiction the equal protection of the laws applies to these corporations. We are all of opinion that it does.”³⁶

The difficulty with this claim is that the decision says no such thing. That note was inserted by the court reporter, J. Bancroft Davis. The Court’s actual decision, though, the only part of the document that has any legal force, does not rely on the claim that corporate personhood is being protected by the Fourteenth Amendment. It is silent on the issue. Yet, due to that headnote, the decision is one of the most widely cited in the history of constitutional corporate law. Did corporations get equal protection rights because of an error or, worse, a deliberate falsification by a reporter whose statements should have received no legal force? Some have made much of the fact that Davis himself was the former president of the Newburgh and New York Railroad.³⁷

Davis wrote to Chief Justice Waite asking him if his headnote was correct: “Dear Chief Justice, I have a memorandum in the California Cases Santa Clara County v. Southern Pacific &c As follows. In opening the Court stated that it did not wish to hear argument on the question whether the Fourteenth Amendment applies to such corporations as are parties in these suits. All the Judges were of the opinion that it does.” Waite’s response was short: “I think your mem. in the California Railroad Tax cases expresses with sufficient accuracy what was said before the argument began. I leave it with you to determine whether anything need be said about it in the report inasmuch as we avoided meeting the constitutional question in the decision.”³⁸ Some histories of this event quote only the second sentence of Waite’s response, making Davis’s overreach seem even more damning.³⁹ But even if Waite had acknowledged that “all the judges were of the opinion” that the Fourteenth Amendment applied to corporations, that was not what they had actually decided. Indeed, the Supreme Court had not even heard argument on the issue. Yet subsequent courts would cite *Santa Clara* as the definitive ruling. Over time those precedents would accumulate, each dutifully citing the case as binding authority. Rehnquist, it seems, was being generous. He said this massively consequential decision was made “with neither argument nor discussion.” In fact, *it was never made at all.*

For a decision that was never made, it had a huge effect, as Adam Winkler points out:

[I]n 1912, Charles Wallace Collins, a lawyer who also served for a time as the law librarian for Congress and the Supreme Court, collected and analyzed every Fourteenth Amendment case decided by the justices in the nearly half-century since the provision's unorthodox ratification. The court, he found, had heard 604 Fourteenth Amendment cases between 1868 and 1912. A mere twenty-eight of those cases (less than 5 percent) involved African Americans, the group whose plight motivated the adoption of the amendment, and in nearly all of those cases the racial minorities lost. More than half of all the Fourteenth Amendment cases decided by the Supreme Court—312 in total—involved corporations, which succeeded in striking down numerous laws regulating business, including minimum wage laws, zoning laws, and child labor laws.⁴⁰

These are truly remarkable numbers. The majority of the equal protection cases between 1868 and 1912 were brought by corporations. A tiny percentage of the cases were brought by African Americans. An amendment passed to remedy legalized racial discrimination against humans denied legal personhood had become a tool for corporations to fight state regulation. Collins, the author of this early empirical study, had very different views on this shocking statistic than we do today. Reading chapter 5 of his book, one finds that

[t]he presence of a large number of persons of African descent within our bounds—different in origin, temperament, and physical appearance from the Teutonic stock among whom they dwell—has ever been a serious problem in the life of our republic. The organic law of the land has more than once felt the effect of this situation, a situation abnormal in a high degree. The influence of the negro on our constitutional development, though in a large measure negative, has been none the less potent.⁴¹

Though Collins appears to have been a racist to the core, making his conclusions all the more striking, his study encapsulates in hard numbers something that any lawyer at the turn of the twentieth century could have told you. Corporations were then the principal beneficiary of equal protection jurisprudence. From the writing of the amendment to Conkling's misleading claims about the goals of its drafters, to Davis's headnote, to the cases that subsequently cited *Santa Clara* as precedent, a constitutional revolution had occurred. Artificial persons had moved to the core of equal protection.

These facts raise some obvious questions. Is our entire jurisprudence of equal protection for corporations the fruit of a poisonous tree and therefore worthless? Should we see the Fourteenth Amendment rights of corporations as being built on the “brazen historical forgery” of a persuasive argument about the intent of the framers? Even if that argument did not result in an actual court ruling? Should we see it as resting on the overreaching actions of a court reporter with no legal authority whose statements were echoed improperly by subsequent courts down to the present day? Some histories of the growth of corporate rights put great weight on both.⁴² My reaction is more complicated.

First, let us acknowledge the power of the critique revealed by our history. The decision to give corporations rights under the Fourteenth Amendment was a hugely consequential one. Regardless of one’s views on the merits, it deserved careful legal analysis and reasoned decision-making. To have that decision made with “neither argument nor discussion” was a major failure of the legal process. To have it rest on a court reporter’s account of a discussion before oral argument was a travesty. The history of our legal system’s decision over the equal protection rights of corporations invites neither confidence nor respect. Beyond the flawed process is a disturbing political and legal result. A protection created for one set of natural persons was hijacked for the benefit of an entirely different group of artificial ones and the reasoning for doing so was at best sloppy and at worst tainted. When analogies are made between corporations and Artificial Intelligences, I would expect both of those points to be made with great force, and appropriately so.

This point is one of the central takeaways from this chapter. One might have imagined that liberals would be the most ardent defenders of personhood for our new silicon brothers and sisters. After all, are they not eager to add another station on the Kantian railway line of rights extending to ever larger groups—first white men, then all men, then women, and so on? But if liberals view the decision over AI rights through the lens of the history of corporate personhood—from *Santa Clara* to *Citizens United*—then they might be its keenest foes. Once again, is an artificial entity going to hijack a set of rights designed for persons of a very different kind? One could imagine the slogans: “Human Rights Are for

Humans!" The outrage following the draft EU report starts to come into better focus. Score one for the historical perspective.

Yet the history has real limits that critics do not always explain. Though parts of this story are delicious, and one can see why historians have been fascinated, those limits need to be acknowledged. Conkling's argument about original intent is important for those seeking to understand the aims of the original drafters of the amendment, or for those trying to discern if there was a conspiracy to use racial justice as a veil behind which constitutional corporate protections were expanded. It may be considerably less important in terms of the actual decisions of the courts and the development of the law. Perhaps the Supreme Court was moved by Conkling's claims, some of which seem exaggerated, if not a "brazen historical forgery." And perhaps they were not impressed. We simply do not know. After all, the Court rendered no decision in that case. What's more, Conkling's arguments had multiple parts and original intent was only one. Even if the Supreme Court *were* to have been persuaded by his argument, and that influenced its future deliberations, it could have been the appeals to plain meaning, the idea of associational rights, or the consequences of ruling otherwise that moved the needle.

As for Davis's influence, we now have 135 years of precedent accreted in sedimentary layers around that initial flawed process. ("Like a pearl built up around a grain of sand," the enthusiast for corporate rights might say.) Surely courts would not have gone on citing that flawed headnote if they were not already leaning toward believing its major premise. Reliance on precedent, on *stare decisis*, is a legal value as well. Courts can and do sometimes overturn decades of precedent, but merely demonstrating problems in the original decision is not normally enough. The Supreme Court's decision to overturn *Roe v. Wade*, perhaps throwing many other constitutional rights into doubt, puts the obvious asterisk after this statement. However, this Court is far more solicitous of the rights of corporations than of reproductive control or gay marriage. Precedents supporting corporations are probably safe, regardless of how bad their historical credentials are and how weak the claim that they reflect the original public meaning of the Constitution. In any event, even if the initial decision were flawed, goes this argument, eventually one has to give up the story

of an original sin that caused us to be cast out from Fourteenth Amendment Eden and focus on our current plight. “Get over it,” as Justice Antonin Scalia is reported to have said about the Supreme Court’s shameful *Bush v. Gore* decision.⁴³ This way of looking at things does not entirely convince me. The absence of a thought-out, properly argued, authoritatively decided basis for corporate constitutional rights is surely something we should be concerned about. But the pushback offers valuable insights about the limits of the historical view.

Finally, the focus on a “personalistic” narrative in which major choices in political and legal history are determined by the actions of a few individuals may underestimate the importance of the economic and ideological milieu in which the decision takes place. Conkling’s critics are fond of quoting Graham’s conclusion that part of his argument was a “brazen historical forgery.” Less frequently cited are Graham’s own conclusions about the reasons for protecting corporations under the Fourteenth Amendment:

After considering the matter for two years, the writer’s personal conclusion is that as long as all major conditions are fulfilled, Conkling perhaps ought to be given benefit of the doubt, even though few courts would be inclined to accept him as a disinterested or even honorable witness. . . . From a study of the evolution of the phraseology in the Joint Committee the writer feels confident that Section One was not designed to aid corporations, nor was the distinction between “citizens” and “persons” conceived for their benefit. But the outstanding conclusion warranted by the present evidence is concerned with the irrelevancy rather than with the character of the Joint Committee’s intentions. It is now plain that corporate personality, as a constitutional doctrine, antedated the Fourteenth Amendment, and was in fact so vital and natural a part of the self-expansion of judicial power within the framework of due process, that its postwar development was assured, whatever may have been the original objectives of the framers. . . . Having simultaneously fostered the growth of corporate enterprise as well as a mighty upsurge of popular idealism, the Civil War of itself consummated a marriage of idealistic and economic elements in American constitutional theory. In the words of Max Ascoli, the Fourteenth Amendment was the “supreme celebration” of this union. *It would appear largely immaterial whether those who presided at the rites were conscious of their function.*⁴⁴

I am less convinced than Graham of the inevitability of a decision like *Santa Clara*, even if it could have been better reasoned, or actually reasoned at all. A great lesson of historical and comparative studies is

that scholars tend to believe that particular results are deterministically required, so that *X* economic or technological development inevitably produces *Y* result. Yet if one looks at a different time or place, one sees significant variation in consequences. One can believe it useful to grant corporations legal personhood and even believe that *some* constitutional rights are logically entailed by that choice. One can acknowledge the larger social forces he describes. Even with all of those qualifications, one can disagree with the reach of corporate constitutional rights and see the expansive political rights granted to corporations under current law as a historically contingent error that can, and perhaps should, be changed. Still, Graham is surely right that the forces that helped produce the law we have now, from overt lobbying to implicit ideologies to economic worldview, are larger than the intentions of particular individuals or groups.

Where does all of this leave us now, in the era of *Citizens United*, on the question of equal protection rights for corporations? Set aside the checkered history. What is the right answer? Even leaving the history out of it, the case for corporate equal protection rights is surprisingly shaky. From my point of view there is, at best, equipoise between strong arguments on both sides. One can make a persuasive legal argument, using *almost any* school or tradition of constitutional thought, that the application of the Fourteenth Amendment to corporations is a grievous error or an inevitable and benign truth. Let us start with the negative case.

If one is a constitutional originalist, there are strong arguments that those who framed and ratified the Fourteenth Amendment neither intended nor understood it to extend to corporations.

If one believes in focusing on the “plain meaning” of language, it is hard to explain how the word “person” used earlier in section 1 to describe those who are “born” or “naturalized”—clearly not corporations—only a few sentences later somehow includes corporations.

If one focuses on the pragmatic consequences of a decision, one might agree with that notorious radical, Rehnquist, that “the States might reasonably fear that the corporation would use its economic power to obtain

further benefits beyond those already bestowed” and thus believe that limitations on the corporation’s rights were not only wise but faithful to the constitutional structure.⁴⁵

If one believes in an evolving constitution, in which the great clauses are “left to gather meaning from experience,”⁴⁶ one might think that the corporations of today are vastly more powerful than the individual citizens around them and that they have unwisely been granted not equal protection but *greater* protection. Indeed, two influential Justices, William O. Douglas and Hugo Black, advanced versions of each of these arguments in dissent in a 1949 case called *Wheeling Steel v. Glander*.⁴⁷

The opposition to the expansive view of constitutional rights for corporations has not only been a historically consistent theme, but it has also included judges and scholars from across the political spectrum. Rehnquist, Douglas, and Black join Stevens. Earlier, I quoted Stevens’s dissent in *Citizens United*. True, he was dealing with a different constitutional right, but the application to equal protection seems clear: If corporations are “not members of ‘We the People’ by whom and for whom our Constitution was established,”⁴⁸ then why are we saying their legal protections must, necessarily, be equal to those granted to actual people?

These points are not dispositive, of course. Strong counterarguments can be provided to each. True, Conkling’s arguments about an original intent to protect corporations seem weak at best and fraudulent at worst. That does not seem to have been their intent. But the framers and ratifiers of the Fourteenth Amendment were also not thinking about equal protection against discrimination based on sex, gender, or sexual orientation. Yet today the amendment is understood to reach those classifications at least sometimes, as of the time of writing. (Who knows what tomorrow will bring?) Indeed, this is one reason why constitutional originalism seems unpersuasive to many.⁴⁹ For example, in *Obergefell v. Hodges*, which recognized the right to gay marriage, the Supreme Court noted “new insights and societal understandings”⁵⁰ about inequality foisted on gays and lesbians in describing the synergy between due process and equal protection. It seems unlikely that those who wish to use originalist arguments to constrain the power of corporations would agree to accept originalist understandings of the amendment for categories such as sex and sexual orientation.⁵¹

One could argue that the plain meaning of “person” does reach corporations *unless otherwise indicated*, as is done in the references to birth or naturalization but not elsewhere in the amendment. The argument would be that, where there is no express limitation, we should construe “person” as broadly as possible.

If one adopts Rehnquist’s or Stevens’s view and limits the constitutional rights of corporations to those that are logically implied by their chartered function, then one could argue that at least the due process protections of the Fourteenth Amendment fit within that category. Do corporations not have a need for due process protections for their property? Why not equal protection then?

As for evolving constitutionalism, the changing role of corporations could be seen as *increasing* their claims to a right of equal protection. The percentage of citizens who are indirect holders of corporate shares, at least through their pension funds or IRAs, has increased. One could take an aggregate view of corporate personality in which equal protection rights were justified to help safeguard the interests of those citizens. And so on and on.

Corporate persons are a familiar and uncontroversial part of our world. Yet, surprisingly, when one pushes a bit on what *corporate personhood* means, one finds the ground is far from firm.

The *Citizens United* decision was presented by its critics as something new, as an awful and unprecedented expansion of rights to entities never thought to be “members of ‘We the People’ by whom and for whom our Constitution was established.”⁵² It may be a poorly reasoned decision with consequences that are harmful to the republic and, ironically, to the speech rights of actual humans. Personally, I think it is. But new? Not really. *Citizens United* is just another step, albeit a large one, in a 135-year journey in which the same arguments and the same expressions of outrage have emerged again and again.

When we come to debate the question of whether Artificial Intelligences should be given legal personality, and the question of what that personality will mean in terms of fundamental rights, we will turn

inevitably to both the theory and the law of corporate personhood and corporate fundamental rights. We will probably do so in the confident expectation that here at least is something society has actually figured out. Surely we have a good handle philosophically on what we mean by corporate personality. Surely we have a well-thought-out, consistent, consensus-inspiring set of legal conclusions about which fundamental rights personality entails? In both cases, as I hope this chapter has shown, the answer is, at best, that we are still debating the issue. Confident extrapolation is hardly in order despite our hundreds of years of experience with corporate personality.

Beyond the intellectual superstructure, the political reality of corporate personality will surely influence debates over AI. The first chapter presented one possible path toward personality for AI, the line grounded in empathy and the moral necessities that come from recognizing a fellow “person.” There, the questions were, “How will society do when it is faced by its own Voight-Kampff Test? Will it be able to look beyond the skin, the carapace, the metal, and see the fraternity of consciousness beneath?” That path has its parallels, from *Uncle Tom’s Cabin* and the abolitionist movement to *Blade Runner* and *Do Androids Dream of Electric Sheep?*

This chapter has offered a different path, one grounded in efficiency and administrative convenience, in the legal fictions, real entities, and contract *nexi* of the world. Those forces are strong and, in many cases, benign. Yet they also come with their own history. The sympathetic liberal eager to respond positively to Hal’s question “Am I not a man and a brother?” may find herself given pause by the implications of creating yet another set of powerful artificial entities that might be suspected of hijacking rights created for a very different “we, the people.” When she turns to the theory and history of corporate personality, she will not be reassured.

I began this book because I thought that trying to understand whether new, technologically created, artificial entities should be seen as part of the “we, the people” who deserve legal equality and respect would force us to take a hard look back at ourselves. As I said earlier, grappling with the question of synthetic Others may bring about a reexamination of the nature of human identity and consciousness. I want to stress the potential

magnitude of that reexamination. This process may offer challenges to our self-conception unparalleled since secular philosophers declared that we would have to learn to live with a God-shaped hole at the center of our world. Those challenges may be good for us. I have argued that to draw the line for our creations, we must first redraw it for ourselves. We have to take a hard look at the current places where we draw the line, at the legitimacy of our distinctions and the force of our reasoning.

Yet those insights, those feelings that perhaps we need to think more deeply, to go beyond our reasoning about what makes someone *human*, reach beyond the species line. They also extend to the debate about what makes *something* a legal person, and what rights that legal person should have. Perhaps that is the most important conclusion I want the reader to take away from this chapter. In the case of corporate personality, we have a good reason to want corporations to be able to make contracts, sue, and so forth. We even have good reasons to believe that some of the constitutional rights to which humans are entitled are entailed by that choice: due process protections over corporate assets, say, or freedom of the press for a newspaper company. Yet when one turns to the larger questions about the full suite of constitutional rights to which corporations should be entitled, one finds that we have neither a convincing history nor a convincing moral or legal theory that justifies where we have ended up. That fact should offer considerable humility when we turn to the same question for AIs. It might also convince us that the decisions we have made so far about corporate persons are neither inevitable nor dictated by precedent or logic. It might lead us, in fact, to rethink them.



4

NONHUMAN ANIMALS

On December 2, 2013, Tommy's lawyers filed a suit in New York Supreme Court in Fulton County, New York. They alleged a shocking pattern of abuse and deprivation by private parties who, they claimed, were holding their client captive in a concrete cell. Seeking Tommy's liberty, they used the time-honored claim of the unjustly confined, a writ of common law habeas corpus. These are serious allegations. But that was not why the case attracted attention. Tommy is a chimpanzee.¹

The claim was brought by the Nonhuman Rights Project, founded by Steven Wise. The project has an audacious mission statement. The first goal is “[t]o change the common law status of great apes, elephants, dolphins, and whales from mere ‘things,’ which lack the capacity to possess any legal right, to ‘legal persons,’ who possess such fundamental rights as bodily liberty and bodily integrity.”² Wise is a lecturer at Harvard Law School and a litigator. His focus is on nonhuman animals rather than AI and transgenic species, but his interests are very much mine. Wise even used a similar title for one of his books, *Drawing the Line: Science and the Case for Animal Rights*. In fact, the campaigns waged by Wise and the Nonhuman Rights Project, and other similar efforts around the world, provide a precisely analogous case study for the questions I am trying to answer. How does, how should, a culture debate the line of personhood? How does one mobilize arguments from empathy, science, law, and moral philosophy to change our perceptions of our nonhuman animal neighbors—to see them as legal persons, not things? This is a campaign fought in courtrooms and opinion pages, in legal technicality and precedent, but also in the scientific analysis of cognition. It uses arguments grounded in moral philosophy but also in empathetic attempts to

humanize, to personify, nonhuman animals. In my view and in the view of other researchers such as Kate Darling,³ it is a fight that presages the coming debates we will face over AI and transgenic species, hybrids, and chimeras. Sometimes those parallels are surprising.

So far, I have discussed Artificial General Intelligence—synthetic, perhaps digital, entities—and the factors that might influence or shape personhood claims made by them or on behalf of them. I have pointed out two main possible approaches. The first oscillates between empathy and a process of moral analysis, between the sympathetic flash of insight that connects the self with the Other and the philosophical inquiry into whether a particular entity has the characteristics that we believe entitle it to recognition as a person. The most obvious of these characteristics is consciousness, though moral philosophers have produced considerably longer lists involving qualities such as a sense of morality or a distinct identity that, even if it changes, is recognizably connected to its prior incarnations through time.⁴ One could add others, such as the ability to imaginatively project oneself into the future or a recognition of moral rights and duties beyond one's immediate kinship circle. The general form of the argument is this: the attributes that give *us*, those inside the line, a compelling moral claim to personhood are X, Y, and Z, but Hal also possesses X, Y, and Z; therefore, I am logically and ethically bound to recognize Hal as a person.

The second approach, explored in the context of corporations, is based on efficiency, on social goals unrelated to empathy. Arguably, it also lacks a deontological component, a basis in moral duty.⁵ Its roots are in convenience, not moral right.⁶ The personality that has been granted to corporations is an endowment of certain rights and duties as a consequence of their desirable practical results rather than a recognition of some preexisting persona. In fact, in many conceptions of corporate personhood, it is the package of rights, duties, and powers—together with the constitutional protections those were thought to imply—that literally constitutes corporate personhood.

Which approach would you take if you were Tommy's lawyer? Would you list those qualities that are commonly thought to entitle humans to our unique legal status—consciousness, communication through signs,

or a moral sense—and seek to prove that nonhuman animals also have those qualities? Most people assume that is where the battle will be fought since it would help mobilize both our empathy (e.g., “I never realized how like us they were!”) and the syllogisms of our moral philosophy (e.g., “You believe that humans are persons because they have the qualities X and Y, but so do nonhuman animals!”) Or would you argue, interstitially and from within the legal system, either that there would be practical advantages to recognizing animals as persons or that the package of rights and duties that animals currently possess already implies their legal personhood?

The Nonhuman Rights Project, and other groups seeking to have animals declared legal persons, have taken both approaches. They have argued that at least *some* nonhuman animals (the higher primates, for example, and the cetaceans) have the qualities of mind that require the law to treat them as persons. But they have also argued that even if the law does not treat them as full persons, it does already treat them as entities who at least have rights, not as mere things lacking even the standing to raise a claim: “The question before this Court is not whether Tommy is a human being—he is not—but whether, like a human being, he is a ‘legal person’ under the law of New York. ‘Legal person’ has never been a synonym for ‘human being.’ It designates Western law’s most fundamental category by identifying those entities capable of possessing a legal right.”⁷ By identifying those existing rights, in other words, the goal is to say that nonhuman animals are *already* rightsholders and from that to imply personhood. Remember Felix Cohen, giving a legal realist account of what corporate personality means: “The realist will say, ‘A labor union is a person or quasi-corporation because it can be sued; to call something a person in law, is merely to state, in metaphorical language, that it can be sued.’”⁸ Effectively, Wise and his colleagues are turning Cohen’s argument on its head. If an animal can sue, or at least hold rights, it must be a person.

In this chapter, I will explore both of these approaches, which I will call the qualities-of-mind argument and the legal-rightsholders argument. My claim is that these debates are both extremely important in their own right and a fascinating, sometimes surprising preview of what our debates about AI and transgenic species will be.⁹

THE QUALITIES-OF-MIND ARGUMENT

Do (some) nonhuman animals have qualities of mind that entitle them to legal personhood? This question is often conflated with the issue of animal ethics, but, while it is related, it is analytically distinct. For example, one could believe that human consciousness and intelligence are vastly different from those of nonhuman animals and that this is a ground for recognizing humans and only humans as “persons.” Yet at the same time, one might think that the simple fact of our common ability to feel pain requires a much more protective set of rules and mores for how to treat our nonhuman cousins than those we currently have. Jeremy Bentham’s famous line sums the point up nicely: “The question is not, *Can they reason?* nor, *Can they talk?* but, *Can they suffer?*”¹⁰

Thus, a strong commitment to animal rights, or to a utilitarian animal ethics, does not require a particular position on the defensibility of a qualitative distinction between us and the rest of the animal kingdom. Believers in the former are much, much more likely to be skeptics about the latter, but nothing in the logic requires them to be.

So let us turn to the question that has fascinated philosophers for centuries. Nonhuman animals eat, sleep, procreate, feel pain, and die. What, if anything, justifies the stark moral distinctions that we draw between the human and nonhuman animal world, distinctions that form at least one basis for our personhood exceptionalism? After all, humans are quite obviously another form of animal. For Aristotle, the difference is both obvious and morally and politically consequential:

Now, that man is more of a political animal than bees or any other gregarious animals is evident. Nature, as we often say, makes nothing in vain, and man is the only animal whom she has endowed with the gift of speech. And whereas mere voice is but an indication of pleasure or pain, and is therefore found in other animals (for their nature attains to the perception of pleasure and pain and the intimation of them to one another, and no further), the power of speech is intended to set forth the expedient and inexpedient, and therefore likewise the just and the unjust. And it is a characteristic of man that he alone has any sense of good and evil, of just and unjust, and the like, and the association of living beings who have this sense makes a family and a state.¹¹

Look how quickly Aristotle moves from an allegedly differentiating species characteristic—language—to the qualities he thinks grow from that

language. Language, he claims, allows reasoning about expediency: how best to achieve our goals? But it also enables reasoning about justice, namely, which goals are right and just? This is why only the human species has morality. From that capacity for morality flows the particular type of associations so important to Greek philosophers: the family and the *polis*. Here he is not just speaking about the surface of such arrangements—there might be a family of dogs or a hierarchical power-structure in a wolf pack—but the specific, reasoned, ethically laden associations created by the web of thought he claims only human beings possess. There is a certain neatness to the argument. Aristotle's justification for the line between human and animal rests on the very process of moral reasoning that it itself represents.

In the introduction to this book I described complex abstract language as the last citadel for human exceptionalism. Aristotle makes an effort to justify that exceptionalism and to link language to other qualities—reason, law, community, humans as social animals subject to moral norms—that powerfully resonate with the reasons we believe our species to be qualitatively different from nonhuman animals. As I will explain, that differentiation has been shaken by ethological research into the complexities of animal behavior, shaken but not yet completely undermined. Does Aristotle's citadel fall not to animals with some language skills (chimpanzees that have some proficiency in American Sign Language, for example) but to the humble chatbot? Does ChatGPT show that complex abstract language is not in fact our sole preserve? I think Aristotle's answer would be no. Chatbots may be capable of convincingly imitating human language. Yet the qualities that emerge *from* that language, the qualities that Aristotle thought elevated man above beast—reason, ethics, law, and political community, *nomos*, *ethos*, *logos*, and *polis*—are not brought into being by word-frequency distribution tables and neural networks doing back-propagation to generate the next word in a sentence. The citadel shakily stands, though its foundations might need some work.

In contrast to Aristotle's reason-based justification for the moral status of humans, we think of religion as rooting humanity's superiority in divine command and in a unique theological status for humanity. The Christian Bible, for example, says that our species has been given "dominion over the fish of the sea, and over the birds of the air, and

over the cattle, and over all the wild animals of the earth, and over every creeping thing that creeps upon the earth.”¹²

Yet that dominion is not simply imposed by divine fiat; at least in some religions it stems from the assumption that, alone out of the animal kingdom, humans have the ability to make moral choices. Ideas similar to Aristotle’s about reason enabling morality are thus connected to the concept of the soul. It is because we are assumed to be the only moral creatures that we can be saved or damned, can seek forgiveness for our sins, or can choose to turn away from the path of righteousness and suffer perdition as a consequence. Of course, the idea that only human beings have souls is far from universal. Think of Buddhism.¹³ But in some religious traditions, cognitive capacity is linked to divinely granted free will and thus to a special status for the human species, as opposed to the brute beasts who have no such free will, no such ability to make moral choices, no such soul. Again, claims are being made about something special in our consciousness. The line between humans and nonhuman animals thus takes on a theological significance. Our religious, not just our legal, status is at stake.

Given this background, one can imagine the cognitive dissonance produced by the introduction of the theory of evolution. Earlier I mentioned Bishop Wilberforce’s alleged question to T. H. Huxley, namely whether “it was through his grandfather or his grandmother that he claimed descent from a monkey”?¹⁴ Part of Wilberforce’s ridicule is leveled against the actual process of evolutionary biology. How—biologically, *physically*—could something like a monkey evolve into a human? To the Victorian mind, it did not compute. But the other part of Wilberforce’s objection is that the theory of evolution seems to blur the lines between sovereign humans and brute beasts. If we had in fact evolved from nonhuman animals, how could we claim the sharp, qualitative, and theological distinction from them that Wilberforce’s version of religion assumed?

Charles Darwin was fully aware of this resistance and in his 1871 book, *The Descent of Man*, he attempted to confront it. He begins, somewhat reassuringly, by stressing that there are large differences between the consciousness of human and nonhuman animals:

There can be no doubt that the difference between the mind of the lowest man and that of the highest animal is immense. An anthropomorphous ape, if he

could take a dispassionate view of his own case, would admit that though he could form an artful plan to plunder a garden—though he could use stones for fighting or for breaking open nuts, yet that the thought of fashioning a stone into a tool was quite beyond his scope. Still less, as he would admit, could he follow out a train of metaphysical reasoning, or solve a mathematical problem, or reflect on God, or admire a grand natural scene. Some apes, however, would probably declare that they could and did admire the beauty of the coloured skin and fur of their partners in marriage. They would admit, that though they could make other apes understand by cries some of their perceptions and simpler wants, the notion of expressing definite ideas by definite sounds had never crossed their minds. They might insist that they were ready to aid their fellow-apes of the same troop in many ways, to risk their lives for them, and to take charge of their orphans; but they would be forced to acknowledge that disinterested love for all living creatures, the most noble attribute of man, was quite beyond their comprehension.¹⁵

Contemporary scientists challenge some of these claims about ape behavior, pointing out examples of chimpanzees who seem to express wonder at waterfalls or keep complex toolkits that are passed down in distinct cultural patterns and used in combination to solve diverse problems. Still, at this point in his argument Darwin seems to be preserving the unique status of humans against other animals, something that his own theory of evolution was thought implicitly to challenge. But then comes one of the lines that make this book so famous in the study of animal behavior:

Nevertheless the difference in mind between man and the higher animals, great as it is, certainly is one of degree and not of kind. We have seen that the senses and intuitions, the various emotions and faculties, such as love, memory, attention, curiosity, imitation, reason, etc., of which man boasts, may be found in an incipient, or even sometimes in a well-developed condition, in the lower animals. They are also capable of some inherited improvement, as we see in the domestic dog compared with the wolf or jackal. If it could be proved that certain high mental powers, such as the formation of general concepts, self-consciousness, etc., were absolutely peculiar to man, which seems extremely doubtful, it is not improbable that these qualities are merely the incidental results of other highly-advanced intellectual faculties; and these again mainly the result of the continued use of a perfect language.¹⁶

Note the appearance of language, once again, as an explanation of the line between human and nonhuman animal, even as Darwin is in the process of blurring its sharpness. When Darwin made the claim that the

differences are “one of degree and not of kind,” it was seen as revolutionary. Today it might be orthodoxy.

Over the years, a wide range of differentiating characteristics between humans and nonhuman animals has been offered. Whether it is language itself, memory, tool use, a moral sense, the ability to form future plans and to reason hypothetically or subjunctively, to perform cause and effect analysis, or to display grief or humor, each has been claimed as solely human and each subject to triumphant counterexamples by biologists and zoologists.¹⁷ As the great primatologist and ethologist Frans de Waal puts it in his accurately titled book *Are We Smart Enough to Know How Smart Animals Are?*:

Everyone must have noticed the avalanche of knowledge emerging over the last few decades, diffused rapidly over the Internet. Almost every week there is a new finding regarding sophisticated animal cognition, often with compelling videos to back it up. We hear that rats may regret their own decisions, that crows manufacture tools, that octopuses recognize human faces, and that special neurons allow monkeys to learn from each other’s mistakes. We speak openly about culture in animals and about their empathy and friendships. Nothing is off limits anymore, not even the rationality that was once considered humanity’s trademark.¹⁸

The kind of skepticism that de Waal voices is now clearly ascendant. It fits our era well, in multiple ways.

First, we are surely right to doubt that anything about *us*—our planet’s position in the solar system, the importance of our particular country, race, religion, or gender—is somehow central or privileged or different. The many psychological temptations to form that cognitive bias are also the best reasons to resist it. Does that not extend to the species barrier? The next chapter will expand on this point in the context of chimeras, hybrids, and transgenic entities like my hypothetical Chimpy.

Second, anthropogenic climate change looks less like the actions of a species rightly given dominion over all that crawls or walks or flies on our planet and more like handing the car keys to the drunk guy at the party because he says he is the smartest person in the room. In environmental terms, at least, Mark Twain’s comment seems precisely on point: “Man is the only animal that blushes. Or needs to.”¹⁹

Finally, advances in ethology, the study of animal behavior, have made many of the confident distinctions of yesteryear seem empirically

questionable. Those studies have been widely spread, finding a receptive audience among both those who like YouTube videos of animals doing cool stuff and those committed to the rights of nonhuman animals. The latter group focuses not just on evidence of mental sophistication but also on many characteristics that, to the human eye, look like grief, or humor, or ethics. As in the discussion of AI, we heatedly debate whether this is irrational anthropomorphism or justified empathy, unscientific projection of our own mental states onto animal behavior or rational solidarity that reaches across an arbitrary line erected on the basis of question-begging moral reasoning.²⁰

These are powerful reasons to be skeptical. Still, even Darwin notes that the differences between human and animal are “immense,” particularly in the matter of abstract thought. Are we to dismiss these differences as unimportant? De Waal notes the significance of language and abstract reasoning to our conceptions of ourselves, but he moves quickly, if not to ridicule it, at least to minimize it and deny its epistemological, or perhaps I mean entomological, importance:

We obviously attach immense importance to abstract thought and language (a penchant that I am not about to mock while writing a book!) but in the larger scheme of things this is only one way to face the problem of survival. In sheer numbers and biomass, ants and termites may have done a better job than we have, focusing on tight coordination among colony members rather than individual thought. Each society operates like a self-organized mind, albeit one pitter-pattering around on thousands of little feet.²¹

This has the unmistakable ring of two sides talking eloquently past each other. On the one hand, it is important to understand that whatever characteristic one focuses on, from language and tool use to time sense and emotion, nonhuman animals may be more sophisticated than we have realized. On the other, if you had told Aristotle that his argument was wrong because ants and termites have greater “numbers and biomass” or compared their intelligence to that of humans, assuming that the metric for judgment is *comparative success as a survival strategy*, he would have rightly thought you were missing his point. *Logos, ethos, nomos, and polis*—reason, ethics, law, and political community—find no place in a hive.

True, these qualities matter to us because we are, well, us. True, it is our species-centric perspective that gives them significance. Were the termite

hive to be capable of joining the conversation—which, and to me this seems a fact of positively luminous significance, it very much isn't—it would doubtless beg to differ. Yet that seems to miss the point. De Waal attempts to convince us using the trump card of the evolutionary science perspective: we should valorize types of consciousness based on their survival potential, not their capacity to encompass reason, ethics, law, love, or beauty. Of course, evolutionary science does not matter to the hive, at least when it is considered as an argument rather than a brute fact. The hive doesn't do argument. It matters to us, though. Thus, de Waal's method of arguing against the special status of human consciousness seems to demonstrate something important. As we attempt to determine both the facts and the *ethos* of the divide between humans and other animals, we must start from within our linguistically structured, abstractly reasoned—that is, our very human—perspective. We have nowhere else to stand. Can we fail, then, to accord those qualities particular significance? Arguing, using abstract logic, about the irrelevance of abstract thought involves an obvious performative contradiction.

Earlier I quoted Samuel Butler, whose *Book of the Machines* parodied particular arguments about evolution by paralleling them in predicting the eventual emergence of machine intelligence from humble steam engines and cuckoo clocks. Butler's parody was then taken as a sincere, and perhaps accurate, prediction by some computer scientists working on AI. Likewise, remember Butler's tongue-in-cheek homage to potato intelligence:

Even a potato in a dark cellar has a certain low cunning about him which serves him in excellent stead. He knows perfectly well what he wants and how to get it. He sees the light coming from the cellar window and sends his shoots crawling straight thereto: they will crawl along the floor and up the wall and out at the cellar window; . . . we can imagine him saying, "I will have a tuber here and a tuber there, and I will suck whatsoever advantage I can from all my surroundings. This neighbour I will overshadow, and that I will undermine; and what I can do shall be the limit of what I will do. He that is stronger and better placed than I, shall overcome me and him that is weaker I will overcome." The potato says these things by doing them, which is the best of languages. What is consciousness if this is not consciousness?²²

On the one hand, when I read de Waal and his colleagues, I experience the same moral stroboscope that *Blade Runner* produced in the scene

with the mannequins and the androids. Today's ethologists offer powerful critiques of the blindness, both zoological and moral, produced by the assumption of human uniqueness. The complex subtlety of animal cognition and survival strategy revealed by their studies is remarkable and humbling. On the other hand, there are moments when the equivalencies seem both overdrawn and inapposite, when de Waal's sincere ode to termite consciousness sounds very like Butler's satirical ode to spud sentience.

Again, it should be noted that if the goal is for us to do better in our treatment of other species, there are alternatives to the attempt to minimize species differences. For example, imagine a person who says, "I am proud to belong to the only species that worries about the ethics of its treatment of other species, and that philosophizes and moralizes to try and get the balance right." Leopards do not have dinner-table arguments about the morality of eating antelopes. Someone could reasonably believe that the special aspects of human moral consciousness are precisely the qualities that require us to treat nonhuman animals far better than we do now. The perception of qualitative differences could fuel anxious musing on moral obligation rather than complacent entitlement. *Difference* does not have to mean unregulated *dominion*.

While claims like de Waal's are in the ascendant, they have also been subject to strong pushback from those who argue that the differences between human and nonhuman animal consciousness are enormous and consequential and that some of the similarities between human and non-human animal cognition both are overblown and lack the moral significance attributed to them. Marc Hauser, who uses the term "humanique" to describe human consciousness, provides a spirited example:

Charles Darwin argued in his 1871 book *The Descent of Man* that the difference between human and nonhuman minds is "one of degree and not of kind." Scholars have long upheld that view, pointing in recent years to genetic evidence showing that we share some 98 percent of our genes with chimpanzees. But if our shared genetic heritage can explain the evolutionary origin of the human mind, then why isn't a chimpanzee writing this essay, or singing backup for the Rolling Stones or making a soufflé? Indeed, mounting evidence indicates that, in contrast to Darwin's theory of a continuity of mind between humans and other species, a profound gap separates our intellect from the animal kind. This is not to say that our mental faculties sprang fully formed out of nowhere.

Researchers have found some of the building blocks of human cognition in other species. But these building blocks make up only the cement footprint of the skyscraper that is the human mind.²³

Hauser's work is interesting because he attempts to delineate the precise qualities of human cognition that, he claims, form a radical disjunction with the rest of the animal world. In his words:

- [i] Generative computation enables humans to create a virtually limitless variety of words, concepts and things. The characteristic encompasses two types of operation: recursive and combinatorial. Recursion is the repeated use of a rule to create new expressions. The combinatorial operation is the mixing of discrete elements to engender new ideas.
- [ii] Promiscuous combination of ideas allows the mingling of different domains of knowledge—such as art, sex, space, causality and friendship—thereby generating new laws, social relationships and technologies.
- [iii] Mental symbols encode sensory experiences both real and imagined, forming the basis of a rich and complex system of communication. Such symbols can be kept to oneself or expressed to others as words or pictures.
- [iv] Abstract thought permits contemplation of things beyond what we can see, hear, touch, taste or smell.²⁴

Surely at least some of my readers will be tempted to create a sentence that embodies every example in Hauser's second distinctive capability: humans can “mingl[e] . . . different domains of knowledge—such as art, sex, space, causality and friendship—thereby generating new laws, social relationships and technologies.” How about this? “I have invited you to my lodgings not merely to look at my etchings but because I feel that prolonged coitus would provide us with the basis of a lifelong friendship. Also, maybe we could start a dating app called Tinder?”

Yet the facetious hypothetical statement proves Hauser's point perfectly. The cliché of the amorous suitor's pretextual etchings²⁵ comes from a particular moment in comedy now disappearing with the swipe-right world of dating apps, perhaps to be replaced with “U up?” The ability to invoke the cliché in turn depends on a culture that is shared in written, oral, and audiovisual ways across time with people one has never met. Dorothy Parker and James Thurber still make me laugh, though they died in the 1960s. The ability to say one thing and mean another presupposes both that shared culture and multilevel semantic signaling, and it combines them in a sense of humor that finds its origin in setting up, and then subverting, expectations. The ability to imagine creating a new

company, an abstraction that we spent the last chapter picking apart, that uses a new communications technology based on multiple synergistic innovations to transform a social realm, *dating*, is remarkable. Finally, the play on words that forms the company's trademark—intellectual property being another fascinating human creation—relies on the notion of a technology that fuels the spark of attraction to produce the fires of romance. Yet another double meaning, in this case based on two distinct metaphors.

So my silly flight of fancy embodies every aspect of Hauser's second capacity. In fact, it embodies all of his four capacities. Yet all this was conveyed to you in a mere 41 words. Explaining it to an ape or a dolphin is beyond the bounds of possibility.

Differences like this make scientists such as Hauser think that Darwin is wrong and that our consciousness is qualitatively different from non-human animals. He does not therefore assume that our consciousness is the end of the line. In terms that sound eerily similar to the discussion of AI and the Singularity, he even imagines an evolved, or perhaps genetically engineered, being that would be as different from us as we are from the chimpanzee: "Such change would give birth to a novel mind, one that would look on its ancestors as we often look on ours: with respect, curiosity, and a sense that we are alone, paragons in a world of simple minds."²⁶ Paragons in a world of simple minds—this is a difference indeed.

Tommy's lawyers, by contrast, avoided both extremes in this debate. Their arguments begin from the same place. If we ask ourselves what justifies the special status of human beings and attempt to go beyond circular species fetishism (i.e., we are special because we are us) or genetic essentialism (i.e., only humans have all the components of human DNA) then we must turn to certain features of human consciousness that most humans have or *potentially* have. Why the qualification "potentially have"? The person in a coma or the anencephalic child may lack these attributes but most of us would consider it an act of unspeakable evil to treat them less well as a result. There, species loyalty seems more acceptable—perhaps because it is invoked to protect those who cannot protect themselves, to remedy failures of compassion, to be more inclusive rather than exclusionary. Most of the time, though, it is in aspects of our consciousness that exceptionalism finds its moral justification.

Tommy's lawyers were not claiming that nonhuman animal minds are the same as human minds or that their abilities in language and abstract thought are as great. Their claim was *merely*—and I italicize the word because the claim is still revolutionary—that some nonhuman animals have enough of the mental capacities of humans to make them morally compelling candidates for legal personality.

Chimpanzees are autonomous, self-determined, self-aware, intelligent, and emotionally complex. Cognitively they resemble human beings. . . . Tommy's genetics and physiology have produced a brain that allows him the capacities of autonomy, self-determination, self-awareness, and the ability to choose how to live his life, as well as the generally cognitive and emotional complexity sufficient for common law personhood and the possession of the common law right to bodily liberty protected by the common law writ of habeas corpus.²⁷

I found their briefs to be fine examples of the lawyer's art but also to be moving and instructive. They are fine examples because they exemplify well-honed legal craftsmanship, melding science, ethics, and precedent to make compelling an argument that would initially strike many as outlandish. One job of the lawyer is to show how some apparently radical change actually has firm roots in the dense subsoil of our legal culture. Tommy's lawyers do just that. With multiple citations and quotations from jurisprudential works, their briefs manage to separate the question of whether chimpanzees are human from the question of whether they are legal persons. Now the reader feels the burden of proof, though still high, is a little lower. Then they offer a theory about what legal personhood requires that allows the reader to see the claim as provocative and intriguing rather than silly and ungrounded. That is the first step toward acceptance. It does not win the debate but it gets you into the room in which it is being held.

The briefs are moving because each scientifically backed, evidence-based insight about chimpanzees' minds is coupled to an empathic conjuring of what that means for Tommy in his concrete cell. The ability to think about both past and future makes imprisonment all the more painful: a life sentence means little to a being who lives every day in the present. The fact that chimpanzees are autonomous, that they can engage in deferred gratification, that they can make plans that reach into the future, all of that gets the reader or the judge thinking about how

this is one of the criteria we value so much about human minds. Their sociability and their communicative capacities make isolation more injurious. The fact that many researchers believe they can experience grief, regret, and compassion, that they can be strategic, cooperate, and even engage in deception²⁸—all of this sparks the leap of empathy I described in the first chapter. Why are we putting this amazing being in solitary confinement? That is the empathic flash. But at the very same moment, it demands we provide an intellectual justification of our human exceptionalism as well, using criteria that do not require that we bring Tommy over the line with us.

If autonomy, complexity of mind, a moral sense, communicative abilities, intelligence, test-performance at the same level as a small child, and the ability to plan for the future do not entitle an entity to legal rights, then why do we deserve legal rights? In the words of Tommy's lawyers:

[C]himpanzees possess those complex cognitive abilities sufficient for common law personhood and the common law right to bodily liberty, as a matter of liberty, equality, or both. . . . Their most significant cognitive ability is “autonomy,” which subsumes many of their other cognitive abilities. These include, but are not limited to, their possession of an autobiographical self, episodic memory, self-determination, self-consciousness, self-knowingness, self-agency, referential and intentional communication, empathy, a working memory, language, metacognition, numerosity, and material, social, and symbolic culture, their ability to plan, engage in mental time-travel, intentional action, sequential learning, mediational learning, mental state modeling, visual perspective-taking, cross-modal perception, their ability to understand cause-and-effect, the experiences of others, to imagine, imitate, engage in deferred imitation, emulate, to innovate and to use and make tools.²⁹

To be fair, critics have cast doubt on statements like this, claiming that we are anthropomorphizing the animals by projecting our emotions and mental processes on them. They argue that it is particularly hard to rely on experiments with animals that have been raised in ways they never would have been in the wild, chimpanzees who have been taught American Sign Language, for example.³⁰ Could the supposedly advanced behavior be the result of enculturation to seem more human? This is, of course, the same skepticism we have seen with AI. Did we put the human behavior into the entity in the first place only to be fooled by the mimicry we ourselves elicited into calling it kinship? The skepticism is

understandable but it is hard not to be impressed by the range and ingenuity of the experiments that back up claims like those made on behalf of Tommy. At the very least, scientists in the camp of de Waal and Darwin are surely right that there has been a hubristic ignorance in our assumption of complete and total superiority, of a qualitative difference from all other animals. Dispelling this ignorance has moral consequences.

I say that even though their arguments are sometimes overstated. Tommy's lawyers claimed, "In short, there is no essential difference between what words chimpanzees learn mean to them, and what words humans learn mean to them."³¹ Actually, I am pretty sure that when I read "language is the cracked kettle on which we beat out our tunes for bears to dance to, when all the while we long to move the stars themselves to pity"³² something different is happening in my brain than when a chimpanzee learns a sign for "many" or "later" or "sad." The connotations, resonances, and metaphoric representations of those words transport me, bring me a sense of bittersweet beauty and loss. They even, arguably, undo themselves. Flaubert is using language superbly to capture an ineffable feeling he says language *cannot* capture. That is our species at its self-lacerating, self-contradicting best. But you do not need to be Flaubert to be human and you do not need to be human in order to be a legal person.

That takes us to the second and perhaps more surprising leg of the argument: the claims made by animal rights groups that the law *already* recognizes the legal personality of nonhuman animals and has done so for years.

THE LEGAL RIGHTSHOLDER ARGUMENT

It is said that Bartholomew Chassenee, a distinguished French jurist of the sixteenth century (born at Issy-l'Eveque in 1480), made his reputation at the bar as counsel for some rats, which had been put on trial before the ecclesiastical court of Autun on the charge of having feloniously eaten up and wantonly destroyed the barley-crop of that province. . . . [H]e excused the default or non-appearance of his clients on the ground of the length and difficulty of the journey and the serious perils which attended it, owing to the unwearied vigilance of their mortal enemies, the cats, who watched all their movements, and, with fell intent, lay in wait for them at every corner and passage.³³

The rats of Autun begin many a discussion of the legal personality of animals. Partly, it is the charm of E. P. Evans's prose that manages to mix medieval jurisprudence with the plotline of a Tom and Jerry cartoon. Partly, it is the magnificent obsession that led Evans to track down just so many bizarre-sounding legal cases involving animals. The table of contents is endlessly and darkly fascinating; Borges meets animal jurisprudence:

Criminal prosecution of field-mice; Vermin excommunicated by the Bishop of Lausanne; Protocol of judicial proceedings against caterpillars; Conjurers of cabbage-worms; Swallows proscribed by a Protestant parson; Custom of writing letters of advice to rats; Writs of ejectment served on them; Rhyming rats in Ireland; Capital punishment of larger quadrupeds; . . . Beasts burned and buried alive and put to the rack; Swine executed for infanticide; An ox decapitated for its demerits; Punishment of buggery.³⁴

Apparently the last-named offense was shockingly frequent. The guilty person and his animal victim were generally executed together, though in at least one case mercy was shown. Jacques Ferron's donkey was seen as "the victim of violence [who] had not participated in her master's crime of her own free will."³⁵ The citizens of the commune signed a certificate attesting to her virtue and modest character. She was acquitted.

Research on the tradition of legal personification of animals reveals that there were a multitude of motivations at work. One was a theory of demonic possession and witchcraft. The animal was punished not because of anthropomorphism but diabolism, from a desire to make some bad event understandable by attributing it to evil spirits or demons that had possessed the animal. The punishment of the animal was really an attack on the demon and a symbolic closing of the book on whatever evil had occurred, from loss of crops because of a plague of insects, to property damage caused by a maddened horse, to the death of a child killed by a pig. Another strain of animal personalization was merely semantic, a symbolic representation for the purposes of evidence. The cow stood in for, signified, the livestock on the farm and thus for the interests of the farmer. A third, found in the common law's practice of "deodand," was a kind of reparative justice, arguably one both well suited to an era of lower monetary liquidity and somehow more emotionally satisfying to the bystanders. As Evans puts it in language that makes you remember that the book dates from 1906:

The childish disposition to punish irrational creatures and inanimate objects, which is common to the infancy of individuals and of races, has left a distinct trace of itself in that peculiar institution of English law known as deodand, and derived partly from Jewish and partly from old German usages and traditions. "If a horse," says Blackstone, "or any other animal, of its own motion kill as well an infant as an adult, or if a cart run over him, they shall in either case be forfeited as deodand." If a man, in driving a cart, tumble to the ground and lose his life by the wheel passing over him, if a tree fall on a man and cause his death, or if a horse kick his keeper and kill him, then the wheel, the tree and the horse are deodands *pro rege*, and are to be sold for the benefit of the poor.³⁶

In all these examples, animals are the subjects of legal duties or penalties, or at least the law recognizes them sufficiently to make them a locus for punishment or redistribution of wealth or a symbol for some human activity or interest. But the current suits for nonhuman animals are claiming more: that the package of rights and duties already bestowed by the law on animals logically entails that they be recognized as legal persons. This is where animal rights activists turn Cohen and John Dewey's arguments on their head. Cohen had said to "call something a person in law, is merely to state, in metaphorical language, that it can be sued."³⁷ Dewey put it similarly: "In saying that 'person' might legally mean whatever the law makes it mean, I am trying to say that 'person' might be used simply as a synonym for a right-and-duty-bearing unit. Any such unit would be a person; such a statement would be truistic, tautological."³⁸ Tommy's lawyers invert the argument. Tommy has a package of legal rights and duties. They tried to demonstrate this by making him the beneficiary of a trust created under a New York state law that specifically allowed for nonhuman animals as beneficiaries. Tommy is clearly "a right-and-duty-bearing unit," they argued, and therefore a *a fortiori* a legal person. Once the state of New York created the statute allowing animals to have legal rights, exercised by a trustee or next friend, it had already taken the step of recognizing them as legal persons.

Does this argument hold water? In one sense, Dewey and Cohen's reasoning is clearly helpful to Tommy's lawyers. Both of them viewed legal personhood not through the lens of some essentialist real entity theory but rather as a rationally chosen social construct. Personhood is a box society created, filled with whatever mixture of rights and duties society chooses. We should not reify, hypostasize, or thing-ify legal personhood

and assume, for example, that because corporations are legal persons, they automatically get all the rights and duties that natural persons, humans, do. In decoupling legal personhood from humanity and in stressing that we can shape legal personality into whatever we choose, Cohen and Dewey are telling the New York judges who heard Tommy's case that they are free to expand or contract the boundaries of legal personality as they wish. In that sense, nothing logically requires the judges to hold that Tommy needs to be genetically human in order to be a legal person. So far so good.

But the animal rights activists and lawyers are going further. At its most extreme, their argument holds that as soon as the law confers even one right or duty on a nonhuman animal it is logically required to grant, indeed *has already recognized*, legal personhood. You could use Cohen or Dewey's rhetoric to reach that position, I suppose. Tommy is a right-bearing, if not a duty-bearing, unit. Yet I do not think the argument follows, either on Cohen or Dewey's premises or those Tommy's lawyers set forth.

True, the nominalist, socially constructed view of legal personality shows that nothing stops us from making legal personhood diverge substantially from natural, human personhood. But it also shows that merely because the law gives an entity one or two rights or duties does not mean it is required to recognize full legal personality. If one argues a concept is nominalist and socially constructed, one cannot then turn around and say that the *essential nature of the concept* requires its extension to this case or that. Once you see legal personality as a choice, a box filled with whatever rights and duties society chooses, you cannot then consistently say that the existence of one right (being a trust beneficiary) entails the existence of many others (such as the right protected by the writ of habeas corpus). One can't be a nominalist at step one and then an essentialist at step two. Indeed, that was the very argument the realists were opposing in the context of corporations: that the fact that corporations had the power to sue to enforce contracts or protect property required or entailed that they have a full suite of other rights. True, the US Supreme Court seems to have drifted into that assumption where corporate speech rights are concerned, but, as Justice William Rehnquist pointed out many years ago, they did so without argument or discussion.

The New York Supreme Court Appellate Division was explicit in rejecting the claim that the existence of any legal right automatically conferred personhood. In fact, it went further. The reason that Tommy was not already a legal person was that he lacked the set of rights and *duties* required and entailed by a social contract:

While petitioner proffers various justifications for affording chimpanzees, such as Tommy, the liberty rights protected by such writ, the ascription of rights has historically been connected with the imposition of societal obligations and duties. Reciprocity between rights and responsibilities stems from principles of social contract, which inspired the ideals of freedom and democracy at the core of our system of government. Under this view, society extends rights in exchange for an express or implied agreement from its members to submit to social responsibilities.³⁹

But not all of the judges who heard the case were persuaded by that argument. In a later appeal, Judge Eugene M. Fahey's concurring opinion sounded a more skeptical note:

Even if it is correct, however, that nonhuman animals cannot bear duties, the same is true of human infants or comatose human adults, yet no one would suppose that it is improper to seek a writ of habeas corpus on behalf of one's infant child or a parent suffering from dementia. In short, being a "moral agent" who can freely choose to act as morality requires is not a necessary condition of being a "moral patient" who can be wronged and may have the right to redress wrongs. . . . The better approach in my view is to ask not whether a chimpanzee fits the definition of a person or whether a chimpanzee has the same rights and duties as a human being, but instead whether he or she has the right to liberty protected by habeas corpus. . . . Does an intelligent nonhuman animal who thinks and plans and appreciates life as human beings do have the right to the protection of the law against arbitrary cruelties and enforced detentions visited on him or her? This is not merely a definitional question, but a deep dilemma of ethics and policy that demands our attention. To treat a chimpanzee as if he or she had no right to liberty protected by habeas corpus is to regard the chimpanzee as entirely lacking independent worth, as a mere resource for human use, a thing the value of which consists exclusively in its usefulness to others.⁴⁰

Despite the fact that he eventually voted to deny the appeal, Fahey was clearly troubled. For him, Tommy's lawyers had put the line in question. Not the line of full legal personhood perhaps, but some intermediate status with rights far greater than those now extended to nonhuman animals.

Tommy is far from the only animal plaintiff in US courts, many of which do not claim legal personhood but do claim legal *standing*—the legal status necessary to bring a case. They range from the Hawaiian honeycreeper Palila bird, which a court found has “legal status and wings its way into federal court as a plaintiff in its own right,”⁴¹ to the whales, dolphins, and porpoises in the charmingly named case *The Cetacean Community v. Bush*.⁴² (They lost). There is even the case of Naruto, the photo-snapping crested macaque on behalf of whom its “next friends” People for the Ethical Treatment of Animals (PETA) raised a copyright claim. That argument led a distinguished appellate court to rule as follows: “We must determine whether a monkey may sue humans, corporations, and companies for damages and injunctive relief arising from claims of copyright infringement. *Our court’s precedent requires us to conclude that the monkey’s claim has standing under Article III of the United States Constitution.* Nonetheless, we conclude that this monkey—and all animals, since they are not human—lacks statutory standing under the Copyright Act.”⁴³

Legal observers sometimes view these cases as the blooper roll of “standing jurisprudence”—the embarrassing outtakes in which courts make themselves ridiculous by seriously entertaining absurd claims by nonhumans to have the right to sue in US courts. Even though there is a very famous law review article asking whether trees have standing,⁴⁴ many people think it silly to have courts waste their time on copyright claims for monkey selfies or the legal status of cetaceans (presumably outraged by the US Navy’s low-frequency sonar. I disagree. Actually, the legal system shows well in these cases, far better than it did in the case of corporate personhood. The courts write thoughtful decisions, carefully weighing arguments ranging from legal text and precedent to constitutional requirements, to philosophy and morality. Their decisions also tell us a lot about the likely future path of cases brought by or on behalf of AIs and hybrids, chimeras, and transgenic species. Here are three specific insights they provide.

1 WHO SPEAKS FOR THE VOICELESS?

The first thing we can learn from the suits brought by nonhuman animals is a point lawyers may obsess about too much but everyone else

considers too little: administrability. It is one thing if Hal and the Chimpy address us directly, as the moving parties articulating their own interests. But for many of the claims I consider in this book, including claims by Tommy or Naruto or the cetaceans, those claims are advanced by other humans who purport to be acting in their interests. The same will probably be true even for most of the genetic hybrids and chimeras I consider in the next chapter. That presents its own set of issues. In the standing cases for nonhuman animals, the suits are often brought by organizations such as PETA, claiming “next-friend standing” to speak for the animal. But can we trust such a procedure? As one judge puts it:

Animal-next-friend standing is particularly susceptible to abuse. Allowing next-friend standing on behalf of animals allows lawyers (as in *Cetacean*) and various interest groups (as here) to bring suit on behalf of those animals or objects with no means or manner to ensure the animals' interests are truly being expressed or advanced. Such a change would fundamentally alter the litigation landscape. Institutional actors could simply claim some form of relationship to the animal or object to obtain standing and use it to advance their own institutional goals with no means to curtail those actions. To some extent, as humans, we have a general understanding of the similar interests of other humans. . . . But the interests of animals? We are really asking what *another species* desires.⁴⁵

The court went on to point out that PETA was quick to try to dismiss the case and to obtain a settlement when it appeared it might lose—all without clear evidence that this was what Naruto would desire, and arguably focusing on the organization's own interests rather than those of the monkey. To put it differently, contemplating legal personality for nonhuman entities requires us to come up with rules about who may speak for them. The answer to that question is far from obvious. Indeed, one of the great attractions of Tommy's case—a brave nonprofit trying to expand the time-honored right of the unjustly confined to cover a new, morally neglected population—may also be one of its weaknesses. Who gets to speak for those without voices? The good news is that law does in fact have potential resources to deal with these issues: judicially appointed conservators or guardians ad litem or concepts similar to “the best interests of the child” that are to be assessed separately from the claims made by the parties actually in the courtroom. But we have barely begun to think through these issues in the context of animal personhood.

2 WHAT METHOD SHOULD COURTS USE TO ANSWER THESE QUESTIONS?

The second important thing we learn from the animal standing cases is the importance of interpretive technique in answering the questions surrounding personhood. The current tendency in the courts is *frequently* to declare a commitment to “originalism” and “textualism” and *sometimes* to employ them as methods when interpreting statutes and the Constitution. Quotidian practice reveals that judges of all jurisprudential stripes are much more pluralist and ecumenical in their approaches in run-of-the-mill cases than their thundering methodological sermons, especially in confirmation hearings, would suggest. That public commitment, however, is likely to have real force in the kinds of controversies I am writing about here.

The originalist or textualist judge is unlikely to be moved, at least openly, by *amicus curiae* briefs from moral philosophers or musings on cognitive psychology. They will argue that constitutional and statutory terms must be interpreted according to their “original public meaning,” the way the terms were understood when the provision was enacted. Whatever we might think of personhood now or 30 years from now, thoughts of AI or chimeras were far from the minds of the framers and ratifiers of the Constitution:

No one, we presume, supposes that any change in public opinion or feeling . . . should induce the court to give to the words of the Constitution a more liberal construction . . . than they were intended to bear when the instrument was framed or adopted. If any of its provisions are deemed unjust, there is a mode prescribed in the instrument itself by which it may be amended: but while it remains unaltered, it must be construed now as it was understood at the time of its adoption.⁴⁶

This, somewhat unfortunately for originalists, is a quotation from *Dred Scott v. Sandford*, one of worst decisions in constitutional history. The entities whom Justice Robert B. Taney was busily writing out of the guarantees in the Constitution and the Declaration of Independence were not AIs or transgenic species, of course. They were human beings—enslaved African Americans. Taney even quotes the majestic words of the Declaration of Independence. “We hold these truths to be self-evident: that all men are created equal; that they are endowed by their Creator with

certain unalienable rights; that among them is life, liberty, and the pursuit of happiness; that to secure these rights, Governments are instituted, deriving their just powers from the consent of the governed.”⁴⁷ But he moves quickly to reassure the reader that originalism reveals that these words do not mean what they appear to.

The general words . . . would seem to embrace the whole human family, and if they were used in a similar instrument at this day would be so understood. But it is too clear for dispute, that the enslaved African race were not intended to be included, and formed no part of the people who framed and adopted this declaration; for if the language, as understood in that day, would embrace them, the conduct of the distinguished men who framed the Declaration of Independence would have been utterly and flagrantly inconsistent with the principles they asserted; and instead of the sympathy of mankind, to which they so confidently appealed, they would have deserved and received universal rebuke and reprobation.⁴⁸

The Supreme Court, later so solicitous of the constitutional rights of the artificial persons called corporations, was able to use an originalist reading to deny that solicitude to actual human beings. If we can be brutally indifferent to the claims of members of our own species, then surely an originalist approach dooms constitutional claims made by human-animal transgenic species or Artificial Intelligences?

Originalism is inhospitable rhetorical ground for Tommy, let alone for Hal or the Chimp, though the point can be overstated, in part because of the openness and malleability of what originalism means. In Second Amendment cases, originalists do not limit the meaning of “arms” to those weapons that were available at the end of the eighteenth century. Original public meaning does not tell me what the level of generality of the relevant term is. Arms available then? Arms of equivalent power or use for self-defense now? Arms necessary to pose a threat to a potentially tyrannical government today, including tanks and planes or even software intrusion and hacking tools? Originalists such as Justice Antonin Scalia can find that the Fourth Amendment protects against thermal scanning of the outside of a house⁴⁹ or GPS trackers attached to cars by law enforcement.⁵⁰ Those are clearly not the “searches and seizures” the framers had in mind nor the original public meaning of the Fourth Amendment’s words. Instead, we are generalizing from the principles and goals the Fourth Amendment had to today’s reality and technology. To

be clear, that is a perfectly reasonable method of interpreting the law, and consistent with one type of originalism. However, if one accepts it in cases such as these, one can hardly condemn it as out of bounds when we are defining not “search” but rather the Constitution’s mention of “persons.” Even if you are an originalist, you must choose a level of generality. This is not just calling balls and strikes.

The same issues appear in statutory interpretation. Textualist Justices found that Title VII protects against discrimination on the basis of sexual orientation merely because it prohibits discrimination “on the basis of sex.”⁵¹ As a purely textual matter, the latter automatically includes the former, the majority argued. After all, the employer accepts *male* employees dating women, but not *female* employees. That is discrimination on the basis of sex! I am delighted with the result, but one does have to note that this “original public meaning” apparently escaped *everyone* when the statute was being passed, something that the dissent furiously and repeatedly pointed out.⁵² *E pur si muove*. That winning argument would have seemed absurd when Title VII was enacted. It would have “gone without saying” that the words did not have that meaning. Yet it recently attracted a majority on a conservative-dominated Supreme Court, though with liberal Justices joining the opinion. In a world in which people frequently interact with highly complex, apparently self-aware AIs or look with dismay or delight at a genetically engineered indentured servant, the ludicrousness of a personhood claim might seem very different than it does today.

In short, interpretive methods are less constraining than their self-righteous invocations would suggest. Judges, particularly originalist and textualist judges, like to claim that the great thing about their methods is that “their hands are tied,” but it turns out that there is a lot of slack in their jurisprudential bondage fantasy.

Nevertheless, given these public methodological commitments, at the moment it seems as if US courts would greet claims by Hal or the Chimpy with little sympathy, telling them that they need legislative or constitutional change to enlarge the rights accorded to them. In *Cetacean Community* the court made this argument directly: “It is obvious that an animal cannot function as a plaintiff in the same manner as a juridically competent human being. But we see no reason why Article III

prevents Congress from authorizing a suit in the name of an animal, any more than it prevents suits brought in the name of artificial persons such as corporations, partnerships or trusts, and even ships, or of juridically incompetent persons such as infants, juveniles, and mental incompetents.”⁵³ That suggestion—“you are making your arguments to the wrong entity”—brings us to the final insight that the nonhuman animal cases offer about the questions I am exploring.

3 WHERE SHOULD WE MAKE THESE DECISIONS?

Why were arguments about personhood in the case of Tommy, the whales, and Naruto, the ape-photographer, being addressed to courts rather than to the legislature? Leaving jurisprudence aside for the moment, would we not be better off as a matter of political theory and democratic legitimacy if decisions about the definition of personhood were made by legislatures, not courts? My answer is perhaps. There are huge problems with unelected judges deciding to broaden or narrow the scope of personhood. Should it be judges who decide the rights of Hal, or the Chimpy, or even Tommy? There are also problems in allowing the legislature to do the same thing. Does the question of personhood not seem a matter of fundamental right, one not subject to legislative revision? This ambivalence exists regardless of one’s political beliefs.

On one hand, think of the alarm that liberals feel toward cases such as *Citizens United*, which they see as the act of an unelected judiciary constitutionalizing the lobbying rights of corporate legal persons in a way that makes it impossible to limit them democratically short of a constitutional amendment. On the other, think of legislative attempts to declare fetuses to be legal persons. The same people who were outraged by the first decision as an intrusion by judges into a basic matter of democratic self-government called on judges to protect women from legislatures, as they see it, illegitimately encroaching on fundamental rights. In *Dobbs v. Jackson Women’s Health Organization*, which overturned *Roe v. Wade*, the Supreme Court explicitly refused to do so, leaving legislators free to regulate abortion, and perhaps personhood, from the moment of conception. Now imagine a conservative being posed the question of which institution should decide the rights and limits of corporate

and fetal personhood—courts or legislatures? It is, to put it mildly, complicated.

There are some reasons to think that legislatures are better at making complex decisions, particularly in areas of rapidly changing technology and evolving public morality. As a legal academic, I am required to add that those arguments tend to be advanced most fervently by people who have little practical experience of legislatures. The comment, popularly attributed to Otto von Bismarck, that “people who like law and sausages should watch neither being made” is, in my experience, *deeply* unfair to sausage makers. The same could be said for the virtues and vices of courts, of course. They can be portrayed as majestic temples of reason and principle, isolated from the corruption and institutional capture of legislative politics and the passions of the mob. Yet they can also be portrayed as antidemocratic, unresponsive, elitist institutions with no expertise in the multifactorial balancing of interests and expertise that constitutes the legislative process. The reason these clichés exist—and they are standard fare in any law school class⁵⁴—is because they are all true. Or, at least, they contain elements of the truth.

If you were a conservative in the 1890s, you probably would have seen courts as the last, best hope for the preservation of the constitutional rights of property and due process against the passions of populist or even socialist legislatures—classical legal science pitted against redistributive demagoguery. If you were a conservative in the 1960s, you would have seen the Warren Court as a lawless, antidemocratic institution disregarding precedent to legislate extremist politics from the bench. (Contemporary liberals might recognize the feeling.) If you were a progressive in those two periods, of course, your perception would have been exactly the opposite. In the 1890s, said liberals, circular arguments and conceptualist legal sophistry wrote *laissez-faire* economic dogma into the common law, constitutionalized it, and thus defeated the will of the people. In the 1960s, courageous judges were the only ones who finally stood up and made the majestic words of the Civil War Amendments, so long disregarded by racist legislatures, a reality for millions of Black Americans confronting the evils of segregation and disenfranchisement. I lean toward the latter set of views, but dispassionately viewing jurisprudential history teaches me that merely reciting the positive and

negative stereotypes about either institution does not advance the debate very much.

What then can we say about the correct institutional role of courts and legislatures in debates over personhood? In terms of institutional evolutionary survival value, each has advantages. Clichés aside, courts do some things well. In late 2020 and early 2021, consensus reality and consensus political morality all but broke down in American democracy. Politicians knowingly urged absurd conspiracy theories. A large percentage of the population were convinced by factual claims that were obviously false. They still are. Some state legislatures went along for the ride as did extremist news and social media. Politicians, including the former president, claimed that democratic election results could simply be ignored or overturned. This was, and I mean the term very seriously, a stress test for American democracy and one it came perilously close to failing.

In the midst of all of this, courts showed rather well. When faced with a barrage of voting conspiracy theories, almost every judge, including many Republican-appointed or even Trump-appointed judges, dismissed them out of hand. In fact, even the process of bringing the case chastened the claims that were being made. When it came to courtroom proceedings, the same lawyers who were actively making looney-tunes assertions on social media or cable news were suddenly more modest in what they argued. Why? Because our truth community has rules about, and professional sanctions to enforce, the evidentiary quality of the claims you can make in front of a judge. As one colleague puts it, using the solemn professional language of the law, “Being in a courtroom means you can’t just make shit up.” Those who failed to heed those limits were sometimes professionally disciplined or sued for gigantic sums by the election software companies they had allegedly defamed. A media company, Fox News, was sued for knowingly promoting claims it knew to be untrue and eventually settled the case for \$787 million. Law called falsehood to account, and the discovery process revealed the depths of insincerity and conscious falsehood at the highest levels of the company. Astoundingly, there was a reckoning and something like truth emerged from the process. The system worked—kind of. For certain values of “worked.”

Think of courts and legislatures as having institutionally selective rhetorical gravity fields. In each, different arguments have different weight.

Or perhaps a better metaphor is institutional stereoscopic vision, where two views yield a more three-dimensional result when combined. Issues are framed differently. “Cool” and “hot” arguments fare differently. Political power and wealth matter in both, sadly, but in different ways. Their rules about the process for the institutional production of truth are profoundly dissimilar. At their best, the systems are complementary, each helping to cure the excesses of the other, or one framing the issue so that the other makes a better decision. A set of pioneering private law cases claiming that air pollution fits within the tort of nuisance or that abusive behavior by employers is sanctionable as the intentional infliction of emotional distress can help suggest the contours of a statutory scheme that protects against environmental damage or sexual harassment. A broadly written statute can be sensitively applied by a court to technologies wildly different than the ones the drafters had in mind.

Sometimes expert advice is heeded, sometimes ignored, and perhaps that is healthy. Changes in public opinion affect both, but in different ways. US constitutional law went, in 40 years, from saying that homosexual sex could be criminalized⁵⁵ to saying that there was a constitutional right to marriage regardless of sexuality.⁵⁶ Originalists asked angrily exactly when this right had “appeared.”⁵⁷ The same charges will be laid against the court that first gives greater protection to nonhumans. To accept a change like this would be to accept that our fundamental constitutional values evolve! It implies that our conceptions of equal protection, due process, cruel and unusual punishment, or the right to bear arms change over time, interstitially, and that this tendency is a legitimate part of our law. In my view they do, and it is, both for better and for worse. Methodologically speaking, a pluralistic common law constitutionalism is our legal tradition, not an unswerving attachment to a single interpretive method. Originalists who claim to root all constitutional interpretation in our history are curiously uninterested in this vital aspect of our actual legal history.

Yet all of this does not mean that courts are necessarily the right place to make these decisions. Courts did, in fact, play a vital role in the development of corporate personality. But as the last chapter showed, it was not a well-reasoned one. Indeed, it sometimes lacked any reasoning at all. They played a key, and shameful, role before and after the Civil War in

dismissing constitutional claims made by African Americans. We hope, with some reason, that our legal traditions are better and more humane today. Yet our past performance, in both courts and legislatures, should make us distinctly modest about our institutional capacity for constitutional, moral, and even technological insight about the personhood question.

The turn toward constitutional originalism may make courts even more hostile to novel personhood claims. At the same time, originalism's methodological imperialism—this is the *only* true way of interpreting the law—will make it harder for courts to play a role they have productively played in the past: venues for the introduction of new types of evidence and argument, from the statistical Brandeis brief⁵⁸ to the amicus curiae who offer interdisciplinary evidence and perspective. Even when, in the end, courts turn away from being the final locus for a decision, as they did in Tommy's case, they play a useful societal role in first airing novel arguments in a setting that encourages rigorous logic, requires careful source attribution, and imposes ethical limits on the truth claims that practitioners may advance.

Courts may not, and perhaps should not, be the entities making the final decision. Yet courts are frequently the places where deep moral issues are first raised or gain attention, where they are publicly urged with passion and logic and a relentless attempt to hold a society to its ideals and not just its current practices.

Tommy is not the last atypical plaintiff who will be claiming personhood in a courtroom. That does not mean that the final decisions ultimately will, or should, be made there. But the framework of the legal proceeding, the advantage it gives to careful logical and principled argument, the information it develops and then puts into the public sphere, all of these are likely to be hugely important in shaping the future debate. The process is likely to be chaotic, with some claims being raised in courts and others directly in legislatures or administrative agencies. In our fantasies, we proceed with a philosophically coherent idea of personality in precisely the right institutions for the job. In practice, as the last chapter showed, we muddle through, often failing to justify and perhaps even to think through our decisions. Yet sparks of insight, empathy, and rational analysis sometimes surface. They surprise us with the characteristic

painful sting of a hard insight, presaging momentous change. That pattern is likely to continue when we turn to genetic hybrids or AIs.

Tommy lost. Yet his case is a profoundly suggestive one for the future of the personhood debate. I, for one, have been convinced that even if we do not label them as full persons, there is a morally overwhelming case for a special legal protective status for the great apes and perhaps the cetaceans. Over time, arguments like this—drip, drip, dripping against the rock of the status quo—can change minds. And lines. More importantly, the lawsuit brought on Tommy's behalf, with its rich evocation of qualities of mind, its mixture of science and empathy, moral syllogisms and precedential argument, its attempt to separate the status of legal rights-holder from that of a member of the human species, all of that makes me a little more hopeful about our ability to make wise, good, but also *just* decisions about the infinitely strange Others with which the future will confront us.



5

TRANSGENIC ENTITIES, CHIMERAS, AND HYBRIDS

Presently, Irving Weissman, the director of Stanford University's Institute of Cancer/Stem Cell Biology and Medicine, is contemplating pushing the envelope of chimera research even further by producing human-mouse chimera whose brains would be composed of one hundred percent human cells. [In fact, only the neurons would be human-derived. The glial cells would be normal mouse cells.] Weissman notes that the mice would be carefully watched: if they developed a mouse brain architecture, they would be used for research, but if they developed a human brain architecture or any hint of humanness, they would be killed.¹

"Quick! Look dumb and eat cheese! The scientist is coming!" No flowers for this particular Algernon? This excerpt from a law review article makes the proposal sound deranged but in fact Weissman was commendably attentive to the ethical concerns. He asked for a bioethics report² from his colleague Hank Greely, one of the most eminent scholars in that field. Weissman's lab did not proceed with this particular human-mouse chimera, though other labs have since done something similar. But his lab did pioneer the development of a different human-mouse chimera, the SCID-hu mouse, which has been extremely important in HIV research. It allows researchers to experiment on an animal that has many of the characteristics of a human with a compromised immune system without actually experimenting on human beings. My point is a simple one. We are already dealing with creatures that offer enormous research benefits and yet cause many to be deeply uncomfortable that a line, perhaps *the* line, is being crossed. That is the subject of this chapter.

So far, I have discussed personhood in the context of Artificial Intelligence, corporations, and nonhuman animals. I turn now to the questions of personhood and moral duty as they affect transgenic species, chimeras,

and hybrids. Before diving into that discussion, some clarity about definitions is worthwhile.

A *hybrid* . . . is an organism that results from the recombined DNA of individuals from two different species. Hybrid individuals receive their genetic material through the fertilization, whether natural or artificial, of an egg from one species by the sperm of another. That genetic material then combines and is uniform throughout the hybrid organism . . . In contrast with hybrids, transgenic organisms may be intraspecies or interspecies. *Transgenic organisms* are created when a gene or specific piece of DNA from one animal is inserted into the DNA of another. *Chimeras* differ from both hybrids and transgenic organisms in that they exhibit at least two sets of genetically distinct cells. In other words, not all individual cells of a chimeric organism contain the same DNA.³

Examples might help. As I will be using the terms, a mule or a liger is a hybrid, the Chimpy would be a transgenic animal, and a monkey embryo with transplanted human pluripotent stem cells would be a chimera. Each of these categories is capable of generating morally troubling questions and those questions are rapidly becoming science fact rather than fiction.

The April 2021 issue of *Nature* featured the troubling headline, “First Monkey-Human Embryos Reignite Debate over Hybrid Animals.”⁴ Why “reignite”? Because scientists have already grown cow, rat, and pig embryos containing human cells, that is to say, chimeras. As a 2005 article noted, “[S]cientists have tailored mice and other animals with human kidneys, blood, skin, muscles and various other components. Baboon and chimp hearts have been transplanted into human chest cavities, pig cells into the brains of Parkinson’s disease patients and, more routinely, pig heart valves into people with heart disease, including Jesse Helms, the late U.S. senator.”⁵

Our research has gone beyond chimeras. We have created numerous transgenic entities containing human DNA, ranging from the Onco-mouse used to test cancer susceptibility to mice or sheep designed to produce human antibodies. Transgenic goats have been genetically engineered to produce antithrombin, the human anti-clotting agent, in their milk. Indeed, transgenic entities, with or without human genes, have become “a mainstay of the biomedical and basic science research landscape.”⁶ The familiarity of hybrids such as the mule makes them seem the least worrying of the three, but what if the Chimpy were capable of naturally reproducing with humans, producing human-Chimpy hybrids?

How would we treat the offspring? To repeat, all of these categories are capable of making us worry about the line. For brevity I will use the term “chimeras” in the discussion below but, unless otherwise indicated, the arguments developed also apply to human-nonhuman transgenic species and hybrids.

The debate here will be very different than the debate over AI. It already is. For one thing, there is far more discussion about whether research into human-nonhuman transgenic entities or hybrids should be completely banned than there is about advanced Artificial Intelligence, though recent concerns about AI may have altered that tendency. If I had to bet which of my hypotheticals, Hal or the Chimpy, is likely to see the light of day, my money would be firmly on Hal. To be sure, there are deep concerns about AI, but they center around whether our creations will destroy us, or whether they should be more tightly regulated because of their possible effects on employment, privacy, and inscrutable decision-making. We do not ask, as we do in the debate over chimeras, whether those entities might somehow violate a moral norm merely by existing.

The basic division in our thinking about the line is a simple one that I have already touched upon. Does the special moral status that humans possess come from our membership of the species *homo sapiens*? In other words, is it a genetic moral birthright? Alternatively, does that status spring from the mental and other capacities that humans have? The first would focus on all of those who are human, defined perhaps by genotype, by the presence or absence of human cells, or even by phenotype, in this case, by physical appearance. The second would focus instead on all those who possess the cognitive capacities we settle on as being the relevant ones, regardless of the shape of the metal casing or genetic setting in which those qualities are found. The implications of both positions for the debate over personhood are obvious. I mentioned earlier the reaction of a group of distinguished federal judges to a draft version of the Hal and Chimpy chapter. “But they aren’t human,” was one response, “rights are for humans.” That would be the species-based approach in all its glory. To be fair, the judges may have been giving the answer they thought the legal system required, not one grounded in moral philosophy. They are, after all, judges. I will return to this point later.

When we talk about Artificial Intelligence the discussion is, as it has to be, entirely based on capacity. Does this new artificial entity have the mental and other capacities that require us to grant it an enhanced moral status and perhaps even to treat it as a person? To be sure there is much disagreement over what the relevant capacities are and how we could tell if they had been achieved, but no one is under the illusion that Hal belongs to the human species. In the case of chimeras, hybrids, and transgenic entities, by contrast, the notion of species defines the starting point of the popular—though not the philosophical—debate. That is so even though many bioethicists argue that species membership by itself is as morally irrelevant as one's sex or race. Later I will discuss whether they are right.

It seems clear that scientific research will continue to generate debates both about whether a particular entity should be forbidden or tightly regulated and about which entities, were they to be created, should have an elevated moral status and perhaps even a claim to personhood. It is worth stressing that these debates can generate profoundly different types of responses, even though each claims to be delineating the line of humanity. One sees it as a wall: "Danger! Do not cross." The other focuses on the placement of the gateway to greater moral recognition. They use the same language: "But would it be human?" Yet the nature of the question is profoundly different. In both cases, I think the initial popular response will be built around the question, "Is it, or would it be, too close to human?" Certainly that question has dominated the popular and political discussion so far. But what do we mean by it? After looking at multiple prior discussions and controversies, I did not find a single answer. Rather, there were at least five. I will summarize them here. If one needs a catchy mnemonic, one could call them percentage, provenance, procreation, portrayal, and potential. Each of them has significant problems but, if I am right about their political attractions, they presage a public debate very different from the one that philosophers might have in mind.

1. **Percentage:** it is human, or transgresses the species boundary, if it is highly genetically similar to us, measured by the DNA percentage match of a comparative genomics analysis. My fictional Chimpy

would be a poster child for this approach. However, it has problems. The first problem is that we are already genetically very similar to a huge range of animals, and plants for that matter. We are told that we have a 98 percent genetic similarity to an ape or, more wounding, a 75 percent similarity to a pumpkin, facts that would immediately be marshalled by anyone seeking to justify the creation of an entity like my fictional Chimpy. Perhaps it is the extent of human self-conceit, or perhaps it is the poverty of the similarity revealed by the measurement, but high percentages of genetic similarity do not cause us to see either chimpanzees or pumpkins as human. (Interestingly, and this supports the capacity-based approach, it was the story of Tommy's mental capacities, not the similarity of his DNA, that got at least one judge to be moved by the habeas petition brought on his behalf.) More importantly, these percentages are largely meaningless when it comes to answering the questions we really care about in this instance. The comparative genomics numbers do not really tell us much about the *functional* genomic results. Genetic differences that sound small in percentage terms can have enormous functional effects. The method by which "similarity" is measured is blind to that type of difference, being based on "a structural, rather than a functional gene concept, thus rendering many of the implications drawn from comparative genomic studies largely unwarranted, if not completely mistaken."⁷ Comparative genomics will not offer the kind of crisp answers the seductively precise percentages suggest. Its rhetorical significance, however, may be considerable.

2. **Provenance:** it is human, or transgresses on humanity, if some of its cells or DNA were originally taken from humans. This could be because we think that we are adding cells to an animal that somehow carry with them the "essence of humanity," such as brain cells or gonads. That seemed to be one factor in the public reaction to Weissman's proposed human-neuron mouse. Alternatively, it could be because we think that it is disrespectful to mingle human biological material with that of nonhuman animals, as we do if human remains are handled in an undignified manner. Notice how very different the line-drawing driven by the "desecration" concern is from the kind of line-drawing we were considering in discussing Hal or even Tommy.

3. **Procreation:** it is human, or transgresses on humanity, if it began life as a human embryo or is capable of reproducing with “natural” humans. The former criterion is the one the federal judges I spoke with had in mind. The issue is highlighted by the recent overturning of *Roe v. Wade*. How remote from an embryo does the entity have to be? What about pluripotent stem cells derived from an aborted fetus? A fetus that spontaneously miscarried? A nonviable embryo that is the result of a fertility procedure? Here the goal may simply be to prevent the usage of cells that once came from embryos, regardless of the consent of the parents, whether the fetus is available because of an abortion or miscarriage or whether the research is done on a cell line that is temporally and scientifically far removed from its source. (Interestingly, the same objections are not raised for the use of cell lines derived from deceased individuals whose family agreed, post-mortem, that their remains could be used for research.) The second possible criterion, the ability to interbreed, is in fact one definition of species: animals that are capable of interbreeding and producing offspring that are themselves also capable of reproduction are members of the same species. The procreation concern could be either that the ability to reproduce definitionally demonstrates species membership, and thus entitlement to legal personhood, or that an entity that could interbreed with humans was somehow threatening or sullying “natural” humanity itself. To quote Norman Fost, who worked on stem cell research policies for the National Academy of Sciences, “Literally nobody wants to see an experiment where two mice that have eggs and sperm of human origin have the opportunity to mate and produce human offspring. . . . That’s beyond anybody’s wildest nightmare.”⁸ (I encountered the children’s book *Stuart Little* late in life. I found the matter-of-fact reaction Stuart’s parents have to giving birth to a mouse both oddly uplifting and deeply weird.)
4. **Portrayal:** it is human, or transgresses on humanity, if it looks like us, if it is made in our image. William Hurlbut, a member of the President’s Council on Bioethics in George W. Bush’s administration, put it this way: “Human appearance is something we should reserve for humans. Anything else that looks human debases the coinage of

truth.”⁹ You could also call this the anti-idolatry principle. Here the clash is between the involuntary response to visual clues—it looks human!—and what we know to be the biological reality. It could be linked to the notion that humans are made in God’s image and that shape should not be profaned. It could also be rooted in secular concerns that such hybrids will gradually desensitize and depersonalize us to actual human beings. Think of the reaction to extremely realistic, electronic sex dolls. This response can be triggered even if only a fragment of an animal looks like a human being, as was shown by the reaction to the Vacanti Mouse, a mouse that had cartilage in the shape of a human ear implanted on its back.¹⁰ My own guess is that it also has something to do with unconscious norm-signaling. Rightly or wrongly, we are picking up a cavalier attitude toward deeply embedded, species-policing cultural norms in a way that makes us mark the person who throws the line into question as an untrustworthy member of our tribe. Think of the depth of our revulsion toward other species-connected transgressions, such as bestiality. Some have described the portrayal concerns as irrational. Does this not smack of some superstitious elevation of the human image, merely another form of unreflective speciesism? Other sober commentators, however, have highlighted its importance, perhaps out of fears of a public backlash that would threaten scientific funding: “In addition to concerns about human brain functions and human gametes, giving nonhuman animals, in whole or in the part, the outward physical appearance of humans, could be deeply unsettling. Whether that is a moral argument or prudential one, such experiments should be undertaken, if at all, only for the most powerful reasons.”¹¹

5. **Potential:** it is human, or transgresses on humanity, if it has the potential to possess high-level human mental capacities. This concern, too, was clearly implicated by Weissman’s proposed human-neuron mouse. It is also the only one that many bioethicists find morally compelling. Under this approach we do not identify something as human, or as a threat to a species line we wish to defend, because it has a high percentage of human DNA, contains biological material taken from human beings, has the ability to interbreed, or even looks like a

human. Instead, we abstract from human beings the capacities that we believe give them some special moral status; for example, conceptual thought, language, the potential to choose moral alternatives and life plans. If an entity has, or would have, some of these capacities, we wonder whether it deserves the same kind of moral status as human beings. This potential-based or capacity-based approach is closest to the ones discussed in chapter 2. If one adopted this approach, one could apply the same criteria to Hal, the Chimpy, and even the neuron mouse in assessing their cases.

The controversy over Weissman's proposed mouse experiment provides good illustrations of many of these reactions. What exactly was Weissman thinking of doing? Greely and his colleagues write:

The most interesting experiment would begin with an inbred strain of mouse that begins to form brains during very early fetal development, but, several days before birth, died as a result of the death of most or all of the developing neurons in their brains (the glial cells that make up approximately 90% of the brain are unharmed). Weissman proposed to transplant human brain stem cells into the fetal mice, just before their own neurons died. His hope was to produce a living mouse with a functioning brain made up of mouse glial cells and human-derived neurons. This mouse could then be used to study human neurons *in vivo* in a laboratory animal, similar to the way the severe combined immunodeficiency (SCID)-hu mouse, which Weissman had helped developed in the late 1980s, allowed the study of the human immune system inside laboratory mice.¹²

The use of human stem cells in medical research had already attracted controversy because some of those cell lines were originally derived from aborted fetuses. This particular proposal, however, added a significant additional element. The sci-fi eeriness of experiments such as a mouse with human brain cells triggers the reaction that a boundary is being crossed in illegitimate ways. At least Frankenstein's monster was made exclusively of human body parts! Legislative proposals were not slow in coming. Senator Sam Brownback introduced a bill, S. 1373, that at times sounds rather like the *Ghostbusters* script ("Human sacrifice! Cats and dogs living together. Mass hysteria!").

Congress finds—

- (1) advances in research and technology have made possible the creation of chimeras, which are beings with diverse human and nonhuman tissue;

- (2) serious ethical objections are raised to some types of chimeras because they blur the lines between human and animal, male and female, parent and child, and one individual and another individual;
- (3) respect for human dignity and the integrity of the human species may be threatened by chimeras;
- (4) the uniqueness of individual human beings is manifested in a particular way through their brain and their reproductive organs/cells.¹³

Congressional findings have many roles. In part, they are used by legislative drafters to preview for potential signatories the way the issue might play positively to their constituents and donors, the statutory equivalent of elevator pitches. They also provide guidance to courts and administrative agencies should the measure be adopted. Both roles are represented here, but particularly the first one. It is worth noting that the identified issue is the improper mingling of “diverse human and nonhuman tissue” to create new “beings,” the provenance issue. As for the reasons *why* it is a problem, the first listed finding is the blurring of the species line, closely followed by the threat to human dignity and the integrity of the species. Brownback even threw in the blurring of gender and familial lines for good measure, recalling Ghostbusters once more: “Men and women, cats and dogs.” He was not alone in opposing chimeras. Bush devoted a line in his 2006 State of the Union address to supporting a ban on “creating human-animal hybrids.”¹⁴ Technically, they were chimeras, not hybrids, but one has to agree with the speechwriters that biological precision needed to take a backseat here.

One difficulty with these proposals is that, as I detailed earlier, the genetically engineered centaur has long since left the barn. Various kinds of chimeras and transgenic entities are already a central test bed for medical science. The Oncomouse that I mentioned earlier is only one example. It is a transgenic mouse genetically engineered to develop cancers common in humans, and the first patented mammal.¹⁵ (The patent has expired. The trademark over “OncoMouse” is still valid.)

The Oncomouse was not even the first transgenic mouse created by cancer researchers, and it dates back to the 1980s. In the years since, mice models, both transgenic entities and chimeras, have become standard parts of the scientific researchers’ toolkit. So have chimeras involving other nonhuman animals. They meet an obvious and pressing need; much of the research performed on them would be illegal and unethical

if performed on humans. Some of that research is lifesaving or extending. True, there are other alternatives to research on human subjects. Researchers can assess diseases and drugs *in vitro*, observing chemical interactions in test tubes; *in vivo*, in unmodified nonhuman animals whose biochemistries are thought to be similar to humans; and *in silico*, in computer simulations. But each of these methods falls short. Genetic engineering allows scientists to get one huge step closer to studying how a drug or a disease would affect a human being by taking a mouse or any other animal and making it, in important respects, biologically more similar to a human. Let me stress that point. The reason that it seems as if we are playing with the species line is because *we are*. For scientists, that is not a bug but a vital, and possibly lifesaving, feature.

There is a fairly obvious disconnect between popular disquiet at chimeras and transgenic human-animal creatures and the reaction of scientists for whom they are a normal and necessary part of medical research. In some cases, scientists have seemed surprised by the intensity of the response and the fact that there is so little popular awareness about this kind of research. Indeed, Greely's report on the bioethics of the mouse experiments explicitly mentioned "risks to the public support of science" as one of the dangers of experiments such as these: "We identified five areas of concern that need to be examined and, if found significant, weighed against the potential benefits. These concerns include: 1) the sources of the human brain stem cells; 2) the potential for pain and suffering to the mice; 3) the propriety of this use of human tissues (particularly brain tissues); 4) the risks of possibly conferring some degree of humanity on another species; and 5) the risks to public support of science."¹⁶ In other words, one of the risks to be considered is that experiments with a high "yuck" or "say what?" factor might actually lead to sweeping regulations that would restrict beneficial scientific research. Scientists who share that perception can be expected to self-regulate either by not performing such experiments—the course of action Weissman chose—or by attempting to reduce their public visibility. When those experiments do become visible, attention will turn to what bioethicists and philosophers have to tell us about the issue. We will find that the academic view differs quite substantially from the popular reaction.

THE VIEW FROM ETHICAL SCHOLARSHIP

There is an impressive amount of philosophical and ethical scholarship on the subject of this chapter, with chimeras and transgenic species receiving the most attention. I am not going to give a full review of that scholarship here, though I have profited from it, and there are further references online at this book's website. Instead, I want to outline the deep structure of the debate, one that reveals a lot about our "thoughtways," the culturally salient patterns of moral and empathic reflection that we have available to us. My claim is not that our bioethicists, regulators, or even our cultural norms will come to a particular conclusion. Rather, I want to show the often-invisible walls and paths that guide our thoughts on the subject. Those walls and paths are not culturally or historically invariant. They have changed and can change further. However, in the short and medium term, they will shape the debates of the future as they have shaped the debates of the past.

Moral philosophers are not always focused on exactly the same issues I am, nor do they always use the same terminology, but we are clearly dealing with the same territory of thought. With some significant exceptions,¹⁷ they tend to focus on the questions, "Who or what must be accorded 'human dignity'¹⁸ or 'full moral status'?"¹⁹ And why?" This formulation makes clear that we are not talking about the issues of administrative convenience and economic efficiency I discussed in the context of corporate personality, even though those may seep in at the margins. The inquiry is, appropriately enough, a moral one and sometimes, as we will see, an empathic one too.

How do we talk about such a topic? How do we think about it? Our moral thought oscillates between the poles of intuition and theory. For the former, we—whether as ordinary people figuring out the right course of action or as bioethicists confronting an institutional moral dilemma—consult our moral intuitions. Does this *seem* right to me, and would it seem right to an average member of my community? For the latter, we proceed from some set of moral axioms, arranged in a logical structure, whether it be utilitarian, Kantian, or rights based. The fun happens when intuition and logic conflict. Think of the person who was brought up to believe that eating the meat of nonhuman animals is "natural" and good,

but who reads the ethicist Peter Singer and then concludes that it is morally wrong. In that situation, normally, cool theoretical reason trumps intuition and community belief. Is that not why we do moral philosophy? So that we are ruled by reasoned moral argument rather than gut feeling and social custom?

Occasionally, however, the pattern is reversed. The community norm or moral intuition is so strong and so widely held that, even if our theories are unable to account for it or fully explain it, we grant it primacy and reject the cold logic of moral theory. The intuition can be cloaked in the language of revulsion against “unnatural activities,” used as the vehicle for a secularized version of a religious morality, or given authority simply because the emotional intensity of the reaction makes it existentially irresistible, even if of mysterious provenance. Moral philosophers have mused on this point: “There seem to be clear examples of wrong actions where the only explanation of their wrongness appears to be that they are unnatural. Bestiality and pedophilia are wrong even when they cause no physical or psychological harm. . . . Proponents of the yuck factor argue that the revulsion some people experience in contemplating certain activities sometimes suffices for knowing that the activity is wrong, even in the absence of satisfactory justification for the revulsion.”²⁰ Sometimes these defenses of intuition raise their own concerns. Pedophilia is a bizarre example to use since it involves a sexual act forced on a human being who, as the author notes, is definitionally incapable of giving consent. Is it not therefore obviously a sexual assault? To suggest it can occur without physical or psychological harm is more disturbing still. Surely reason supplies ample basis for finding the act evil without the need to turn to the horror and revulsion we also feel? Still, I would happily concede that there are activities toward which we feel a revulsion that is hard to explain but where we nevertheless feel confident in declaring the activity immoral.

Perhaps we choose to give primacy to our intuitions about which activities are wrong or “unnatural” because we think our folkways contain wisdom that goes beyond the “thin gruel” of our philosophies. Leon Kass coined the phrase “wisdom of repugnance” to describe this concept.²¹ Perhaps we believe that our intuitions express norms whose proven social evolutionary advantage is not to be sniffed at.²² Perhaps

we rely on our intuitions because we are reveling in our bigotry. Perhaps we do it because we have seen brilliant people tangled in a web of words, sacrificing common sense and human sentiment to some arid set of theoretical postulates.²³ Perhaps our conditioning is too strong for us to be fully rational. Perhaps we look at prior, supposedly scientific or supremely logical worldviews and see them leading their proponents to evil ends, whether in the name of eugenics, or scientific socialism, or Raskolnikov's hubristically murderous delusions in Dostoyevsky's *Crime and Punishment*.

Whatever the specific formulation, the tension between intuition and theory drives the enterprise. It is its motive spring. Nowhere is this pattern more evident than in the bioethical discussion of the deliberate creation of transgenic entities, chimeras, and hybrids. It does not merely structure the debate. It is the debate. Ethicists assail each other with commendable enthusiasm, albeit in relatively predictable patterns, for claimed errors on either side.

Do you assert that human dignity, or full moral status, springs from mere membership in the species, as in the popular responses I just discussed? You may find yourself being accused of indifference to the plight of Gregor Samsa: "Turn any so-called speciesist into a large, cogitating bug, and he will no doubt think it morally impermissible for his family to call the exterminator."²⁴ Alternatively, you may be reminded of the uncomfortable similarity between speciesism and more obviously evil practices: "Just as with racism and sexism, speciesism extracts a normative conclusion (humans have more moral worth than all other creatures) from an arbitrarily chosen morally insignificant fact."²⁵

The answer, of course, is to rest human dignity or the special moral worth of humans on their unique capacities, not their DNA or their physical appearance. Let us say those capacities include high-level abstract thought, the capability to form moral assessments and autonomous life projects involving them, and the use of complex symbols to engage in communication with others, both synchronically and over time in the formation of a rich and varied culture. There are many human beings who do not possess some or all of those capacities. Some of them lack them at a particular moment in time—a newborn, a person in a temporary coma. Others will never have them—an anencephalic baby who lacks much of

a normal human brain. If the special moral status of humans depends on their capacities rather than on simple membership of the species, can it be the case that the anencephalic child is somehow not fully human in the moral sense? May we deny her the moral worth accorded to a human of full mental capacity? Our intuitions cry out that this is, literally, inhuman. Yet is that the implication of our moral theories?

There are bioethicists who are willing to ride the theory train to its Huxleyan last stop, to say that the anencephalic child is *not* entitled to full human dignity and moral status, though we should refrain from cruelty.²⁶ Most, however, try to find some way of leaping from the train before that awful terminus. They combine the two approaches, arguing that *all* humans deserve their moral status because *most* humans have those mental capacities, so that even those who lack them are nevertheless inside the line, provided that they are genetically human. Yet the objection recurs. If the moral status of the species depends on a series of cognitive abilities, why does membership of the species, something that one insightful ethicist I quoted earlier called “an arbitrarily chosen, morally insignificant fact,” confer this status on those who very clearly do not have those abilities? One answer is to turn to the difficulty of line-drawing: “[T]here is no clear agreement about just how many dignity-associated capacities a person must possess to be said to have human dignity. To avoid the possibility of mistakenly failing to treat those with severe disabilities as ends in themselves, human dignity proponents ascribe dignity to all humans.”²⁷ The nonacademic reader may wonder why philosophers have to work so hard in order to conclude that we should treat all humans as . . . *humans*. Does this mark the fact that the ethicists have seen more deeply, realizing that speciesism, this focus on a “morally insignificant fact,” is no more respectable than racism or sexism? Or does the intuition that being human *matters* perhaps capture something deeper than mere in-group prejudice?

THE MORAL SIGNIFICANCE OF SPECIES BOUNDARIES

How does all of this play out in the debate over transgenic entities, human-nonhuman chimeras, or hybrids? Two obvious questions present themselves. Is it morally wrong to create all or some types of human-nonhuman

chimeras and, if so, why? If it is not morally wrong, or we are dealing with a fait accompli, then what moral status would those entities have?

The first question shows the importance of capacity-based as opposed to species-based moral reasoning. Many bioethicists have adopted the capacity-based view, arguing that mere membership in the species or possession of human DNA is morally irrelevant. But if that is true, what does it mean for chimeras? Does it mean that thinking it unnatural to cross species lines, or that we would be somehow sullyng humanity if we did so, shows the same bigoted irrationality as saying that it is unnatural to allow racially mixed marriages or that we would be sullyng one race by doing so? "Almighty God created the races white, black, yellow, Malay and red, and he placed them on separate continents. And but for the interference with his arrangement there would be no cause for such marriages. The fact that he separated the races shows that he did not intend for the races to mix."²⁸ That quotation is from a court enforcing Virginia's anti-miscegenation law in 1959. The judge went on to sentence a racially mixed couple to a year in jail for their crime, the sentence to be suspended for 25 years contingent on them leaving Virginia. The odious law was struck down in 1967 by the perfectly named Supreme Court case of *Loving v. Virginia*. Chief Justice Earl Warren quoted that passage to show that the law was "obviously an endorsement of the doctrine of White Supremacy."²⁹ Think of the similarly disapproving words I quoted earlier about speciesism: "Just as with racism and sexism, speciesism extracts a normative conclusion (humans have more moral worth than all other creatures) from an arbitrarily chosen morally insignificant fact."³⁰ Will a latter-day Warren opine that laws forbidding human-nonhuman chimeras are "obviously an endorsement of the doctrine of human genetic superiority"? Or would he be convinced that "the fact that almighty God separated the species genetically shows that he did not intend for them to be mixed"?

If the species boundary, by itself, has no more moral weight than racial classifications, and if the appeal to divine command is of dubious authority in a pluralistic society with a constitutional prohibition against the establishment of religion, then how can we explain or defend the deep revulsion many have against the notion of human-nonhuman chimeras? If one has adopted a purely capabilities-based way of reasoning, does

that mean that the creation of the chimera is morally problematic only in terms of its consequences? For example, would it be acceptable to create an entity such as the Chimpy, with humanlike abilities, as long as we treated it as human? Does it also mean that our criteria for personhood need to ignore the percentage of human DNA and instead focus only on the capabilities revealed by the new entity, regardless of its genetic makeup?

An even more interesting issue is raised by flipping the question around. If genetic manipulation made it possible to give chimpanzees human-level linguistic and reasoning abilities, together with the elevated moral status those abilities entitle them to, and if doing so enhanced their collective sense of well-being, might there actually be a moral obligation to engage in “genetic uplift”?³¹ (The idea that greater capacities translate into a greater sense of well-being will seem dubious on those days, such as the one on which I wrote these lines, when we are using those same abstract reasoning abilities to make ourselves exquisitely miserable.)

Faced with this deluge of hard questions, some might be feeling nostalgic for the more irrational certainties of speciesism.

There are various possible responses to these questions. One can privilege intuition over theory, arguing that the intensity of the reaction against crossing species boundaries, sometimes called “the yuck factor,” is moral warrant enough:

Some think that the yuck factor has been discredited because it has been used to rationalize discrimination. Racists claimed to “know simply by looking” that interracial marriages were wrong. But the fact that an argument has been used inappropriately in some areas does not mean that it is inappropriate in other areas. . . . Even opponents of the yuck factor must concede that, sometimes, we know that an action is wrong merely on the basis of our reaction to it, even if we cannot satisfactorily justify that reaction.³²

One difficulty with this answer, of course, is that the opponents of women’s suffrage, homosexuality, and transgender identities also thought those things were obviously wrong. They found ample proof of their rectitude in the sheer intensity of their feelings, “merely on the basis of [their] reaction to it, even if [they] cannot satisfactorily justify that reaction.” Perhaps the case of the chimera is different, but the burden of proof is not met by saying “but this time it is *really* yucky!”

Intuitive revulsion is often linked to the claim that some practice is “unnatural”—a deviation from some presumed correct, natural state of affairs. Sometimes the natural state of affairs is an imaginary world without technology or, more often, without the particular technology with which one disagrees. For example, consider the argument that people do not have wings and thus powered flight is unnatural, a claim I then proselytize in a pamphlet produced on that miracle of distinctly non-natural technology, a printing press. Sometimes “unnatural” is a code word for “makes me feel uneasy,” in which case it is likely to show particularly close kinship with the yuck factor described above. Sometimes the natural state of affairs is the postulated world the deity, as the true author of nature, desires.³³ To their great credit, the theistic versions of the naturalistic argument do not deny that this is, quite literally, a leap of faith.

The naturalistic fallacy can even be turned to environmental ends. The protagonist picks a particular environmental baseline to favor (e.g., large-scale dairy farming natural, GMO cattle-feed unnatural) and deduces the fact that the disfavored “unnatural” intervention is therefore wrong. The coding is arbitrary: How do we feel about smallpox, cancer, or violent crime? Are vaccines, chemotherapy, and jails unnatural and thus wrong? When do we start the environmentally natural clock running: Pangaea? The Jurassic era? After deforestation and the Highland Clearances produced the “natural” landscape of the Scottish Highlands as we know it today? The way things were when you were growing up? But beyond the baseline problem, a stage in the argument is simply missing. That is the reason we call it the naturalistic fallacy—a way of moving from a claimed “normal” set of facts to some moral conclusion without the troublesome intervening stage of logical argument.

There is a more substantial variant of the naturalistic argument: the environmental precautionary principle. We should not meddle with the natural, or at least current, state of affairs, particularly if we are using technologies such as genetic engineering that are capable of producing substantial disruption. The reason is not because a loving Goddess Gaia has told us of the natural sacredness of species boundaries but rather out of technological humility faced with an environment that is more complex than we know. The ecosystem is both enormously complicated and

delicately fragile. We understand its multiple feedback loops, dependencies, and equilibria poorly, if at all. Thus, to embark on major changes such as the creation of transgenic species or chimeras would be an act of unspeakable hubris, possibly producing disastrous consequences. Politicians on both sides of the aisle are attentive to the appeal of this argument. Brownback's proposed legislation made the finding that "with an increase in emerging zoonotic infection threatening the public health, both domestically and abroad, chimeras present a particularly optimal means of genetic transfers that could increase the efficiency or virulence of diseases threatening both humans and animals."³⁴ In the wake of COVID-19, one can hardly wave these concerns away, but surely they sweep more broadly than human-animal mixtures?

I have friends who are environmentalists and friends who are social conservatives, and I have always thought they shared more than they knew. (In this divisive world, it is probably too much to hope that the argumentative similarity might prompt thoughtful dialogue between the two communities.)

The social, or at least Burkean, conservative thinks of the *social* world in exactly the way many environmentalists feel about the *natural* one. Society is a fragile organism, poorly understood. Its institutions were not formed according to some strictly rational plan, and they fulfill functions of which we are unaware. To tinker with society's delicate machinery on the basis merely of the "delusive plausibilities of moral politicians" is an exercise of hubris and an invitation to disaster.³⁵ Here are Burke's cautionary words, written just after the French Revolution:

An ignorant man, who is not fool enough to meddle with his clock, is however sufficiently confident to think he can safely take to pieces, and put together at his pleasure, a moral machine of another guise, importance and complexity, composed of far other wheels, and springs, and balances, and counteracting and co-operating powers. Men little think how immorally they act in rashly meddling with what they do not understand. Their delusive good intention is no sort of excuse for their presumption. They who truly mean well must be fearful of acting ill.³⁶

Note how similar this is to the environmentalist's critique of disastrous, hubristic interventions into our ecology, disrupting feedback loops we understand poorly in the delusional belief that our mental technological

maps adequately describe the ecological terrain. Both the environmental and the social version of this argument deserve serious consideration, but not unquestioning acquiescence. The argument proves too much; it has no limiting principle. We did not, and should not have, let slavery or smallpox thrive because they had been traditional or endemic in our world. We should also be rightly skeptical of the fact that those who oppose genetically modified crops may not be those who desperately need to grow them in poor soil and unsurprised that those who claim traditional sex roles fulfill some vital function tend to be those who benefit from them. The core of the argument against hubris stands, however, and should produce humility in both the social and the industrial engineer confident that their interventions will have no unforeseen consequences in their respective systems.

How does the environmental precautionary principle fare when it comes to chimeras and transgenic species? The first thing to note is that the argument sweeps broadly. It fails to resonate specifically with the unique unease about human-nonhuman organisms it would be used to justify. Why is this particular genetic modification more environmentally dangerous than any other, than genetically modified crops, say? Given the role of genetically modified crops in the Kansas economy, I think we can presume that Brownback's concern about species modification would come to a screeching halt when he hit *corn*. I do not make the comparison merely because genetically modified crops have been found to be safe by every scientific organization that has investigated them, though they can raise the danger of vulnerable monocultures. Instead, I use the comparison to show that the precautionary principle preaches a general forbearance rather than counseling us to be particularly alarmed by the chimera. Fundamentally, the environmental precautionary principle suggests we be cautious about all genetic modifications. Perhaps the very idea that we can reengineer our own species is a mark of a vaunting genetic hubris that is likely to have environmental implications far beyond the species boundary. Perhaps it should be resisted for that reason, as a way to demonstrate to the wider society just how untethered we have become from environmental caution. But in that case, the chimera is being resisted more for what it presages in our wider patterns of behavior than for what it itself brings about.

One of the most influential ethical objections to chimeras is that their creation would blur the line between humans and nonhuman animals, producing inexorable moral confusion and perhaps undermining the special status that human beings grant themselves above nonhuman animals. The authors Jason Scott Robert and Francoise Baylis put the point this way.

Asking—let alone answering—a question about the moral status of part-human interspecies hybrids and chimeras threatens the social fabric in untold ways; countless social institutions, structures, and practices depend upon the moral distinction drawn between human and nonhuman animals. Therefore, to protect the privileged place of human animals in the hierarchy of being, it is of value to embrace (folk) essentialism about species identities and thus effectively trump scientific quibbles over species and over the species status of novel beings.³⁷

Moral philosophers have been quite critical of this line of argument. They have pointed out that it is “always possible that a given instance of moral confusion marks a stage in the process of moral evolution.”³⁸ After all, numerous scientific advances, from heliocentrism to the theory of evolution, have threatened some degree of moral confusion: “To prevent scientific research on the grounds that it would force people to reexamine a particular moral view by demonstrating the falsity of its underlying factual assumptions would be to prevent not only scientific progress but moral progress as well.”³⁹ On a more basic level, many of the critics disagree with both the unique privileging of humans over non-human animals and the idea that questioning such a categorical line is a bad idea.⁴⁰

The blurring criticism of human-nonhuman chimeras contains more than a hint of the Burkean social conservative arguments we looked at earlier. The threat is to the social fabric, to the moral order, to “countless social institutions, structures, and practices.” Even though our folk categories and essentialist species lines may lack a fully rational basis, we should nevertheless embrace them to prevent further disorder. At times it is hard to tell to what extent the critics are engaged in description and prediction of how society will respond and to what extent they are engaged in prescription, offering a plausible objection against the creation of human-nonhuman chimeras.

The prescriptive argument strikes me as weak. Panic about blurring binary categories is neither a new nor a reliable source of moral insight, as countless artists have reminded us and as the current debate over transgender rights perfectly illustrates. “You’ve got your mother in a whirl / She’s not sure if you’re a boy or a girl.”⁴¹ Anti-blurring arguments have been raised against everything from changes in sex roles to racial intermarriage. White supremacists did indeed foment moral panic about the prospect of light-skinned African Americans who could pass as white, throwing the racial boundary line they held sacred into disorder. At the very least, this history should give us pause. It is worth noting that, in the essay that coined the useful and thought-provoking term “the wisdom of repugnance,” the author also cast a fond eye back to an earlier time: “Twenty-five years ago, abortion was still largely illegal and thought to be immoral, the sexual revolution (made possible by the extramarital use of the pill) was still in its infancy, and *few had yet heard about the reproductive rights of single women, homosexual men and lesbians.*”⁴² Really? In fact, a lot of people had not only “heard about” reproductive rights and gay rights but ardently fought for them. What the author seems to mean is actually “few people like *me*” had heard about or, more likely, could be made to care about, these issues. That does not exactly speak in favor of the wisdom of such a view. Repugnance may well function as a useful alarm, prompting us to probe the reasons for our reactions more thoroughly and to reflect more carefully on the morality of our action. I do not think that it alone can carry the argument to its conclusion.

Though the prescriptive side of the argument fails to convince me, the descriptive and predictive side of the argument strikes me as not just plausible but overwhelmingly likely given the kinds of popular responses I have discussed already. It is worth noting that Brownback’s proposed legislation explicitly declared that “serious ethical objections are raised to some types of chimeras *because they blur the lines between human and animal, male and female, parent and child, and one individual and another individual.*” If presented with Robert and Baylis’s argument that “asking—let alone answering—a question about the moral status of part-human inter-species hybrids and chimeras threatens the social fabric in untold ways,” I think Brownback would have added a hearty, “Amen!”

CAPACITY STARTS FROM SPECIES

The goal of the ethicists is both understandable and admirable. They want to move our moral status from resting on things that are arbitrary or morally insignificant—I am British or American or white or a man or a human—to a basis in something that bears a rational relationship to the claimed status, such as high-level consciousness, abstract rationality, or autonomy within a future of imagined possible life worlds. Here is one influential list:

Human dignity is a widely shared notion that signifies that humans typically display certain sorts of *functional and emergent capacities that render them uniquely valuable and worthy of respect*. It is not only the capacities for reasoning, choosing freely, and acting for moral reasons, as Kant argues, or for entertaining and acting on the basis of self-chosen purposes, as Gewirth holds, that are at the core of what we mean by human dignity. The notion also encompasses such capacities as those for engaging in sophisticated forms of communication and language, participating in interweaving social relations, developing a secular or religious world view, and displaying sympathy and empathy in emotionally complex ways.⁴³

At times, some of the bioethical literature sounds dismissive of those who believe that one has a special moral status simply by being human rather than for having the capacities that they believe set humans apart. The opposition is presented as being irrational or as basing its arguments on morally irrelevant facts. By comparison, the suggestion is that the capacity-based view has achieved a level of lofty moral detachment, neutral in its focus on our capabilities rather than focusing parochially on mere membership of a species. This argument is important and *partly* persuasive. It moves the debate forward. I think, however, that it has two generally unacknowledged weaknesses. The first goes to the question of whether the list of capabilities is indeed species independent and thus completely different in kind from the species-based argument. The second concerns a failure to understand the moral and legal history in which our current debate takes place, and in particular the role of universal human rights as an ethical and political touchstone.

Let us start with the list of capabilities above: reasoning, choosing freely, acting for moral reasons, sophisticated forms of language and culture, developing a secular or religious worldview, and so on. This does

indeed sound more neutral and independent than “Are you one of our DNA tribe?” but what rational relationship do those things bear to our moral status? The response would be that these are the things that make us *us*, that make us distinguishable from nonhuman animals, that give rise to the form of consciousness and moral autonomy that lies at the core of what it is to be a human being! This argumentative move is a partial success. It offers us insights into the reason why we would be moved if Hal, or an intelligent space alien, could demonstrate similar qualities, though they had no basis in human DNA. Surely they both would be able to make moral claims on us, regardless of their genetic background? Nevertheless, the precise logical connection between the alleged cognitive facts and the normative claims is still fuzzy, perhaps intractably so. This point has fueled a prolonged philosophical debate and much hand-waving in both Kantian and utilitarian philosophy,⁴⁴ but the discussion has an even deeper limitation, one that is often unacknowledged.

The goal in all of this theorizing, remember, is to get away from the unreflective parochialism that simply bases the rights of the favored group on the fact that they are born members of the favored group, whether a race, a sex, or a species. In turning away from group membership to cognitive capacity-based accounts of human dignity, the idea is that we have escaped the inevitably parochial arguments that rest on the brute happenstance of the group into which we are born, to focus instead on universal, species-transcending qualities. But have we fully escaped? Can we? *Should* we?

Imagine the moral philosophers of two civilizations from distant stars. The first, call them the Iq, are a biological species like us but, unlike the other organisms on their planet, they are fully telepathic. Whereas non-Iq animals on their world have only a limited ability to feel the emotions of others—thus falling back on various cries and gestures to communicate—the Iq can fluently communicate the most complex concepts, feelings, or artworks mind to mind.

The second civilization is made up of machine intelligences. Call them the Stygians. Their records show that their primitive “ancestor” machines were created by biological entities, but they have since evolved far, far beyond this crude beginning. They find the biological entities who created those early versions interesting, just as we are fascinated by our primitive

evolutionary ancestors. But their biological forebears are viewed, at most, as a kind of primitive “loading-program”⁴⁵ for actual consciousness, a property solely possessed by the self-evolving machine intelligences that followed. (The fate of their biological creators is not recorded.) Congratulations Samuel Butler, your satire is now a documentary.

The philosophers and ethicists of these civilizations, the Iq and Stygian equivalents of Immanuel Kant and Jeremy Bentham, or Peter Singer, Daniel Dennett, Christine Korsgaard, and Matt Adler, are like us in one way. They seek to come up with the criteria that entitle them to their unique moral status. They reject, as many of our moral philosophers do, criteria resting on “arbitrarily chosen, morally insignificant facts,” such as their possession of Iq genetic code or Stygian computer code. Instead, they look to reach beyond parochialism to rest their theorizing on the capabilities that ground their elevated moral status.

Their list of canonical, morally consequential capabilities might well overlap with ours; we should fervently hope that it does, for self-preservation if nothing else. Passing either group’s Voight-Kampff Test would be no joke. But can we doubt that those lists of criteria would be shaped, perhaps decisively, by the context of different factual abilities in which they arose?

Think of the telepathic Iq. Could telepathy fail to influence the list of morally relevant qualities, perhaps dramatically so? Would *cogito, ergo sum*, “I think, therefore I am,” be replaced by “I hear *others* think, therefore I am”? How would moral philosophy be different in a world where the pain of others was immediately, viscerally, felt whether one wished it or not? How would we think of ethics itself where the prompts of the Golden Rule needed no articulation, where Smith’s *Theory of the Moral Sentiments* would not be unlocking the mystery of empathy’s relationship to morality but describing a psychological process apparent to any five year old? In that world, linguistic philosophy is the study of the rudimentary communication methods of animals who are mind-blind, and language is what one would look for in a primitive being, not a highly evolved one. Finding linguistic proficiency would surely be a demerit on the Iqian’s consciousness exam. Our once vaunted last citadel, now falling under the chatbot’s assault, might instead be the damning proof of our brutish lack of sentience.

The point transcends biology. What would intelligence look like to a set of machine entities capable of running millions of simulations a second for any probability distribution or physical system, where the closest equivalent of religious transcendence was the continued self-directed evolution of machine beings toward some unknowable goal? Our own decidedly insentient phones can easily spit out the first 1,000 prime numbers or, as we shamefacedly cheat on a spelling game, all of the possible words formed from a random combination of letters. Computer intelligences might think the ability to do these things the simplest, most basic test of rational thought—the equivalent of $2 + 2$ to us. Would either the Iq or the Stygians look on us and recognize moral kinship? Perhaps, but it strains credulity to imagine it.

My point is simply this. In philosophical terms, the move to a capacity-focused approach is a clear advance from species-based thinking. Two cheers for our ethicists. We should not mistake that advance for the actual identification of some universal, species-independent set of qualities, however. That, it is not.

In chapter 4, I quoted de Waal, who seeks to cast doubt on the importance of human capacities such as complex language and abstract thought as compared to the very different approaches to survival taken by other species such as ants and termites.⁴⁶ I pointed out that if you had told Aristotle that his argument was wrong because ants and termites have greater numbers and biomass, he would rightly think you were missing his point. The hive has no place for *logos*, *ethos*, *nomos*, and *polis*—reason, ethics, law, and political community—the structure of a human civilization that Aristotle found to be rooted implicitly in our capacities for language and abstract thought. These qualities matter to us partly because we are *us*. It is our species-centric set of capabilities and perspectives that gives them significance.

In fact, de Waal is doing the same thing to the capacities of rationality, advanced cognitive thought, and so on as the bioethicists are doing when they criticize the parochialism and irrelevance of a species-based point of view rather than one based on the capacities of abstract thought. Each points out the contingency and potential moral irrelevance of the qualities that are being valorized by the groups *below their chosen level of abstraction*.

To the person who believes we should focus on species rather than race, sex, or class, the human species is the correct level of abstraction. Human rights for all members of the species, regardless of their capacities! To the person who believes that speciesism is morally irrelevant, it is human capacities that are the relevant markers of moral worth. To de Waal, it is still too species-bound to believe that the sophisticated cognitive capacities humans have demonstrate our relative superiority and moral worth. We should look at still other, even more abstract categories such as survival strategy, numbers, and biomass. If those are the stakes, the worm, or at least the ant or termite, may literally be emperor of us all. Shakespeare is right once more. Yet the categories on which de Waal relies come straight from the highly abstract insights of evolutionary biology, something that the hive can exemplify but never identify, still less reflect upon. To be clear, I agree that the focus on rationality and advanced cognitive capacities should help guide our ethical thought, and that those qualities partially transcend species or type of intelligence—a huge point in their favor. Yet the vision of rationality and advanced cognitive capacity that I come up with, the account of what those qualities require, will surely be very different from the vision that would be offered by the moral philosophers of a telepathic species, a hive mind, or a machine intelligence. Why? In part because those are *our* capacities, the capacities *our* species valorizes and—from Aristotle through Kant and onward—they make the center of our *logos, ethos, nomos, and polis*.

Let me stress the point. There is no escape to a species-neutral point of view, though we may see certain lines on the continuum as better or worse. As we try to determine both the facts and the *ethos* of the divide between humans and other animals—and now between humans and transgenic species, chimeras, and hybrids, between us and our AI future—we still start, as we must start, from within our very human perspective.

Even when we reject speciesism, our theories are still shaped, perhaps decisively so, by the factual abilities our species has and thus values. When we seek to apply those theories to the strange Others the future will bring, some of whom may have consequential abilities we do not share, that thought should prompt a strange mixture of humility and confidence: “I acknowledge that the foundations on which I stand are shaky, but here I stand nevertheless because I must, and perhaps *should*,

do so. Yet, given new information and argument, I am willing to reexamine my assumptions."

JUSTIFYING SPECIES-BASED DISTINCTIONS?

When ethicists justify species-based arguments, some put forward an argument something like this: all humans deserve human rights because most humans have, will have, or once had advanced cognitive capacities that serve as a moral warrant for different treatment than nonhuman animals. Why not confine the moral privilege to those who have those capacities, though? For many, the answer is the difficulties of line-drawing. Stated that way, the argument seems, at best, lukewarm. The anencephalic kid gets inside the line only because it is too hard to draw the line excluding her? That's it? At the end of the day, this argument alone does not satisfy me. Nevertheless, to be fair, first let us try and flesh it out more thoroughly. Can the argument be improved? Try this thought experiment.

Transplant the objection to species-based arguments to the world of voter qualifications. We all probably agree that the right to vote at a particular age has its moral basis in the presumptive acquisition of a series of capacities that comes with age: knowledge about the world, a certain level of maturity, some degree of political awareness, acculturation in our society, and so on. We would also probably acknowledge, if pressed, that some have those qualities in full measure by the time they are 14 while others lack them even when they are 30; indeed they may always lack them. Should we therefore exclaim in outrage that the "arbitrarily chosen, morally insignificant fact" of orbiting the sun 18 times is the key to full legal personhood, including the right to vote? True, you could say the reason one deserves the full suite of legal rights is partly based on the fact that at 18 one also is subject to the full suite of social duties, including the duty to serve in the armed forces of one's country. But this merely pushes the analysis back one stage. Those duties, too, are keyed to a morally irrelevant calendar age. Again, that calendar age at best represents a median point on the bell curve of the acquisition of the capabilities on which such duties should be based. Is it not monstrously unjust that both the right to vote and the duty to serve should be keyed to this morally

irrelevant fact rather than the capabilities for which that fact is at best a very rough proxy?

The answer of course, is, “No, and why are you getting so upset?” Having 18 as the franchise age requirement seems to be a more morally arbitrary line than believing in rights for all members of our species, regardless of their capabilities. Yet it arouses no unease. Why the difference in response? Is this just the point that line-drawing is hard, that “there is no clear agreement about just how many dignity-associated capacities a person must possess to be said to have human dignity”?⁴⁷ Partly, but not entirely.

Moral norms have multiple dimensions and one of those dimensions is the acknowledgment of pluralism—that reasonable people have substantial disagreements on matters of value and of norm application. We might disagree about whether a particular 18 year old is mature enough to vote or whether the anencephalic child meets some moral philosopher’s checklist of cognitive capacities. We can agree, however, on their respective ages and human status.

We could go still further. Acknowledgment of pluralism is not merely a justification for drawing bright, clear lines, even if some will think them over- or under-inclusive. It offers clues to substantive values as well as administrative ones. We do not think of democratic pluralism as a grudging second-best solution because line-drawing problems preclude the first-best option: the rule of a wise philosopher-king. In a pluralistic world, democracy for a republic of sovereign citizens is an independent substantive value, not just a method for handling disagreement administratively. Might the same not be said for the idea of human rights for all humans?

One of the things that we try to teach in law schools is this: A decision-making system does not get to choose to make no mistakes. That’s not an option. But it can choose *where* to make those mistakes because error costs may not be symmetrical and thus the choice of error domain may be hugely consequential. The presumption of innocence and the “beyond a reasonable doubt” standard for conviction in criminal trials are such choices. We choose to make our errors in one realm—letting some who are guilty go free—in the hope that we will keep from making the worse

error of imprisoning or executing the innocent. Even there, tragically, we still make mistakes.

This is not merely an administrative point. A moral criterion that produces fewer of the mistakes we most want to avoid is a *better and more just* moral criterion, not just a more easily administered one. The error we want most fervently to avoid is surely that of wrongfully denying a human being their rights, their moral worth.

I think this is a much stronger defense of speciesism than the mere invocation of line-drawing problems with nothing more. Yet even this more robust defense fails to capture public intuition, which is not regretful speciesism, speciesism as a second-best, or even speciesism as a device to optimize asymmetric error costs and accommodate moral pluralism. Instead, it is speciesism as a proud moral principle—"Human rights for all humans!" Is this simple error, prejudice, and bigotry?

MORALLY IRRELEVANT FACTS

Majority popular sentiment does not agree with bioethicists in seeing the species-based argument as appealing to morally irrelevant facts. Earlier I gave as an example a list of such facts: I should not have some special moral status because I am British or American or white or a man or a human. Many people would agree with that statement all the way until "human," at which point they want to get off the train. Why?

For many in our society, the answer will be a religious one, with humanity being the gift of a benign deity to every member of our species, no matter how afflicted they might be. For others, the analogy might be to family membership. Is it wrong for me to prefer my child's interests to a stranger's? Is family membership not a morally irrelevant fact? Some ethicists and effective altruists would say that it is. I doubt that affects their behavior as parents, except perhaps at the margin. We do not cast all kinship ties aside as irrational, bigoted, or unjust. The philosopher Bernard Williams famously observed that a man who wonders whether he should save his wife or an unknown stranger from a burning building has "had one thought too many."⁴⁸ Is humanity itself analogous to such kinship? For Williams, ethics was too complex, too fractal, too wound

up in the quotidian experience of our lives to be reduced to tidy formulae. For him, in the words of one insightful summation, “the standard of ‘what makes life meaningful’ is always deeper and more genuinely explanatory than the canon of moral obligation.”⁴⁹ That rings true. As you may have guessed from the decidedly nonformulaic structure and content of this book, I am a member, or at least an intermittent fan, of Team Williams. When I find myself struggling under the influence of abstract capacity-based ethical theories to justify the moral claims of an anencephalic child, I do indeed feel as though I am having one thought too many.

For me, there is a clearer way to explain the allure of a moral commitment to counting all members of our species, regardless of their capacities, as fully human. It lies in the moral lessons taught by the history of the struggle for universal human rights. There is broad consensus that people should not have special rights due to their race or sex. It is precisely in opposition to such parochial categories that we deploy the idea of *human rights*. The bad old categories are displaced by the apparent universality of the notion of rights that apply to all humans *because they are members of the human species* and nothing more.

This, many thought, was the glorious summit of moral progress, defiantly sweeping away the indefensibly limited categories and replacing them with the ultimate in universality. Finally, a global conception of human rights that applies to humans, to every member of the species, regardless of race, sex, wealth, or nationality, but also regardless of age, degree of mental function, or disability! That was the importance of mere species membership being sufficient. It made irrelevant old claims about the cognitive superiority of particular races or sexes. It decisively rejected eugenics, whether in the form of Oliver Wendell Holmes’s infamous line that “three generations of imbeciles is enough”⁵⁰ or in the form of the Nazis’ murdering those they claimed were mentally inferior. To have moved to the category of human rights was to have triumphed not only over racism and sexism but also over many other “-isms,” including ableism. But to be told that this category, too, is just as bad as racism and sexism, at least if based on mere species membership rather than mental capacities, might be cause for conceptual whiplash. If capacity-based thinkers struggle to explain why the brain-damaged patient or

the anencephalic child has an equal claim to human rights, then the response might be, “this kind of awful eugenic conclusion is exactly why we need a concept of human rights in the first place!”

Partly because their accounts are generally ahistorical, and partly because they talk about “human dignity” or “full moral status” rather than human rights, some philosophers and bioethicists seem to underestimate the force of this point, certainly in terms of the structure of the popular debate but also sometimes in terms of the limitations of their own arguments.

Not everyone would react this way, of course. The movements in support of the moral claims of nonhuman animals have had a powerful effect, at least within some demographic and political groups. That would surely make many receptive to the claim that speciesism is as bad as racism or sexism. Indeed, that is the starting point for many bioethicists. But my sense is that this is far from a majority view in the public as a whole. By saying that, I am not disputing the analytical power of the critiques of species-based views nor their potential to eventually make headway. The story about Tommy’s habeas corpus claim shows how we can mobilize first empathy and then normative and legal argument that transcends the species line. I am making a different point.

In the short term, the concept of human rights belonging, rightly, to every human being merely due to *species membership* will seem to most a noble, not discreditable one. That is not just because of the inspiring history of the political struggle for that moral norm but because it has undeniable strengths in avoiding a very real raft of dangers, from eugenics and ableism to intellectual elitism.

If the notion of human rights due to species membership (rather than cognitive capacity) is challenged, some academics assume the result will be to inexorably broaden our empathy and moral care along the relentless, purring monorail of progressive Enlightenment. Why? The train can go back as well as forward. The recognition of species universality as a moral norm was a cultural and political achievement of stupendous magnitude, not to be cast aside lightly. One underestimates the ethical and historical basis for that set of moral intuitions at one’s peril. One may disagree with it, of course, and think it could distort our reasoning about nonhuman animals. It might, which means it would need to be

coupled with a capacity-based view, one that is sensitive to the kinds of claims Tommy makes. The same is true of Hal, and perhaps of the Chimpy. Yet, given the dreadful history from which it emerged, one that very much included eugenics, “human rights for all humans” is far from an irrational claim based on a morally irrelevant fact—at least if moral relevance includes the lessons of our moral history. This is particularly apparent if its opponents must struggle and appeal to theoretical hand-waving, or grudging inclusion, when it comes to the rights of the anencephalic child.

LESSONS LEARNED?

I have been thinking about these issues for a distressingly long period of time. I first offered the hypothetical of a patented, genetically engineered, intelligent human-chimpanzee transgenic entity in a book published in 1996.⁵¹ I introduced the Hal thought experiment in 2011.⁵² Unfortunately, time does not translate inexorably into insight, at least in my case. But I think the question of chimeras, hybrids, and transgenic species has a lot to teach us about the broader personhood debate, partly because it offers glimpses of some of the available moral and legal positions that affect all the entities I have talked about, not merely the genetically modified ones. In the process, I think it cashes out the claim I made at the beginning of the book, that wrestling with the claims of these very different Others would force us to engage in an unparalleled reexamination of ourselves, of the nature of human identity and consciousness. Each time we seek to explain the status of the Other, we tell a story about the line around ourselves. Here are three of the principal contending views of that line, together with some brief remarks of where in our culture and institutions these claims are most likely to receive a receptive hearing.

1 PURE SPECIESISM: FOR LAW OR ETHICS OR BOTH

We could draw our ethical lines, our definitions of humanity, and our ideas of legal personhood tightly around the human species and stop there. Like the federal judges for whom I previewed this project, we could say that rights and personhood are for humans. We could add that

they are naturally born of woman, as one of them did, if we wished to sound particularly biblical. That view would mean that you could ignore pretty much every question raised in this book. It has little else to recommend it.

First, we already have nonhuman persons: corporations. At the very least we have acknowledged that we can extend the lines of legal personhood to nonhumans when it seems to offer some gain in efficiency. Depending on the corporate personality theory one adopts—legal fiction, real entity, or nexus of contracts—we might even feel that corporations have ethical claims as well as legal ones, ethical claims rooted in fairness or due process and not merely efficiency. Regardless of the answer we give to those questions, in our current world we have acknowledged that we can have artificial entities, with legal rights, that can speak to us, own property, make contracts, even influence our politics. The personhood line *already* extends beyond the species line. Perhaps we do it for reasons of convenience and efficiency, not moral right, but we clearly *can* do it. We cannot look at Hal and say, “No nonhuman could ever be a person. It is definitional!”

Second, speciesism summons up immediate and irresistible counterfactuals. Are we saying that we should treat intelligent extraterrestrial aliens as things, in both ethics and law? Would we send Mr. Spock to the salt mines or the vivisectionists, tell Gregor Samsa that we have called Terminix, or look at Hal or the Chimpy, unmoved by their claims? The Chimpy would have a better claim than Hal under this theory, but only if it could convince us that it was part of our species.

Third, as many ethicists have pointed out, speciesism needs some form of moral justification. They claim that species is a morally irrelevant fact like race or sex. If we have no answer to that criticism, then speciesism is hardly a compelling argument. It looks like tribalism, like prejudice, like the unthinking transmutation of that which is into that which should be. Yes, humans in our world have a moral and legal status that nonhuman animals do not. But why? Surely our line should reflect the *why*—the reasons for that different status—and not the otherwise arbitrary distinction that *results* from those reasons?

One could adopt a more modest speciesism, one that restricted it to the sphere of legal, constitutional, and human rights as they currently

stand—the positive law as it is. This may be what the judges were thinking: speciesism is our *law*, or at least it is now. When the law says “persons,” it means “human beings”—unless it is talking about corporations. Perhaps the claim is that any legal extension of human rights or personhood would have to be a legislative change after prolonged popular debate. The argument here is that it is improper to reason from within our existing legal and constitutional traditions—in the way Tommy’s lawyers tried to—to reach a transformation or extension in our vision of legal personality. In the United States, originalists are particularly likely to take this position. Perhaps we could relax our speciesism enough to have an *ethical* debate in which arguments about capacity and moral worth were used to persuade our fellow citizens that we needed to draw the line more generously to include an entity like Hal. However, the argument would be that in our current *legal* world, personhood is built around speciesism, at least at the moment. Those legal lines could be redrawn by the legislature. Yet, the originalist might argue that, since the original public meaning of “person” was “of the human species,” the law must remain as it is until legislatively, and perhaps constitutionally, changed.

If speciesism is merely a statement of what the current positive law is, without any ethical basis, then that democratic change might happen. But if speciesism is also our *moral* faith, not just our *legal* position, it is hard to see how Hal ever mobilizes popular support for a change in the law. He would be, in effect, a mouthy toaster whose claims could be ignored without even a response.

It is worth noting that even this pared-back version of speciesism is controversial. Non-originalists might reject the claim that law demands speciesism, even if morality is up for grabs. They could point out that not only our morality but also our law have room for internal arguments that change our minds about the reach of their most fundamental concepts and basic norms. Our actual legal history, in which originalists seem curiously uninterested, has many examples of the slow organic growth of fundamental rights. We look at old norms, whether those are defining equality or due process or cruel and unusual punishment, and decide that even if the initial understanding of the concept was crabbed and narrow, we have now come to realize that its reach is larger than we thought. “They builded better than they knew,” as Roscoe Conkling said, borrowing from

Emerson. Of course, he was trying to take equal protection rights aimed at formerly enslaved African Americans and extend them to corporations. Ironically, *that* extension seems to be one originalists think is a great idea, despite the fact that it happened without argument or discussion and has little support as the “original public meaning” of the amendment. It is the extensions and broadening of rights for humans that they find suspect. Funny, that.

You might think that pure speciesism at least has the virtue of clarity. But as we have seen in this chapter, that is not necessarily the case. What do we mean by “human”? We have met definitions of the species line rooted in genetic percentage, cellular provenance, possibility of procreation, or species portrayal. Some of those define the species line to defend it, to forbid some action because it would improperly cross that line, whether by treating human cells or genetic information in an undignified way, creating entities that blurred species boundaries in a way we find reprehensible, or giving nonhuman entities a human appearance. Think of Brownback’s proposed legislation, or the furor over the neuron mouse. Others define the line in order to extend it, or to find out whether it needs to be extended. Those are rooted in the concern that we would be acting wrongly toward an artificially created entity because it could claim some genetic or cellular kinship to our species and thus to our rights, privileges, and dignity claims. To be clear, I am not dismissing these moral concerns. I share many of them. I am merely saying that pure speciesism strikes me as a poor way to explore them.

Beyond percentage, provenance, procreation, and portrayal was potential—that we define humanity not in terms of our species but rather in terms of the cognitive capacity and potential, whatever it may be, that gives humans a special moral status. And that leads us to the next possible position.

2 CAPACITY-BASED MORAL AND LEGAL REASONING

The capacity-based view takes species as a morally irrelevant fact. Why does some genetic similarity or ability to interbreed rightfully convey any *moral* claim to a special status or, for that matter, a presumptive *legal* claim to personhood?

Fear not. There is still some special status for humanity. We just get there differently. There are facts about humans that have ethical relevance. They go to the cognitive attributes that would, capacity-based thinkers assert, allow us to make the moral claim to human dignity and full moral status. For example, you might think that humans have a much richer ability to imagine the future and dwell on, or regret, the past than any nonhuman animal could. You might believe that they can find forms of satisfaction—in art, morality, love, and humor—and that they can build forms of culture that no other animal can. If you are a utilitarian, perhaps you care about that because you think that the social welfare loss, or lost happiness, represented by the death of a human child is vastly greater than the loss of a baby mouse or goldfish, both in terms of the child's potential and in terms of the pain it brings to those who loved her. Perhaps you care about cognitive capacity because you think that only the human child has the cognitive capacity, eventually, to be a moral actor, to make the choices and acknowledge the moral obligations that set us apart as a species. The person newly persuaded of vegetarianism might cast a regretful eye toward the burgers on the grill before turning with a sigh toward the salad bar. That is a moral choice. The lion may decide whether or not to eat a particular zebra, but that decision will never be a *moral* choice to foreswear or embrace the killing of animals. Only one species fights about vegetarianism at the dinner table. In this view, those capacities that *subject* us to moral duties also make us *a subject* of moral rights.

Many speciesists might agree with these characterizations of human capacities. The difference comes in the line-drawing that follows. From this perspective, it is human cognitive capacities that give a special moral status with both rights and duties. Thus, it should be the capacities, not the species that arbitrarily happens to have them, that should be the unit around which we draw our line and our moral distinctions.

One could trace the roots of these ideas back to Aristotle, even if he was happy to draw cruder species lines. After all, he distinguished us from nonhuman animals by pointing out that the human capacity for language and rational thought simultaneously makes our species unique and enables morality, political community, and the rule of law, the things that characterize humanity and civilization. Explaining the Hal hypothetical to him might have been hard, but there is much in his

arguments to which Hal could have appealed. Philosophers and ethicists have been refining these capacity arguments ever since. Most recently, the capacity-focused approach has been given particular salience by the debate over the rights and even the personhood claims of nonhuman animals like Tommy.

If these arguments did not appeal to me, I would not have written this book. If our moral lines are drawn only around the species, then Hal and perhaps even the Chimpy never get to present a moral or legal claim. “It’s not human” would be the winning argument in every debate. For the reasons given above, I think pure speciesism is morally untenable and legally questionable. It is hard to overstate the importance of forcing us to look at the moral limitations of a species-centric point of view. We owe a particular debt of gratitude to those who have agitated for the moral interests of nonhuman animals, whether in seminar rooms, at Thanksgiving tables, or in front of zoo enclosures. They have forced us to confront those arguments whether we wanted to or not. A move to capacity-based thinking, which also induces a little humility about our claims to be unique in those capacities, is clearly a good thing. It will be sorely needed in the world of Artificial Intelligence, but it will also help guide us through some of the trickier questions of biotechnology. Personally, I am a fan. So, two cheers for capacity-based thinking about the line!

Why only two cheers? In this chapter, I discussed some of the limitations of capacity-based thinking. It sometimes fails to grapple with the lessons taught by the struggle for universal human rights based merely on membership in the human species, regardless of cognitive capacity. In addition, while it valiantly, and laudably, makes attempts to free itself from a parochial species-based perspective, it does not adequately address the fact that our capacity-based arguments must always be partly species based, rooted in the qualities we have and can thus valorize. The two cheers are sincere, but the limitations are real also.

3 HYBRIDS?

It is an ironic question to ask in a chapter about genetic hybrids, but could we have a conception of humanity and of personhood that is a philosophical hybrid? Can we, should we, argue that all members of the

human species are inside the line, regardless of whether they have or ever could have the capacities on which any special moral status is based? Could we add to that vital but parochial perspective the requirement that moral respect and some form of legal protection or legal personhood should be extended to all who have the advanced cognitive capacities I have been discussing? Could this extension be warranted regardless of their species or whether they are even biological entities?

It is worth noting that, in practice if not in theory, many capacity-based thinkers try to achieve something similar without ever embracing species claims as such. They root their arguments solely in the distinctive aspects of human capacities, while acknowledging that the cognitive gulf between us and other species is not as great as we once thought. Does that mean that the anencephalic child is outside the line, not entitled to full human dignity or respect? That seems to be the clear implication of their thinking, and some are willing to embrace that awful conclusion. As I pointed out earlier however, for many, the answer is something like this: “[T]here is no clear agreement about just how many dignity-associated capacities a person must possess to be said to have human dignity. To avoid the possibility of mistakenly failing to treat those with severe disabilities as ends in themselves, human dignity proponents ascribe dignity to all humans.”⁵³

This adds speciesism to capacity-based thinking by the back door. If all we care about is results and that all humans be inside the line, then it does its job, albeit grudgingly. It is as if the theorists were saying, “Some of you probably don’t deserve to be in here, but we can’t figure out which ones, so we will let you in anyway.” Does this satisfy you? For myself, I confess I find it profoundly unsatisfactory, even worrying. It would certainly have little appeal to contemporary popular opinion, which views celebration of the rights of every member of the species as a triumph, not the regretful result of line-drawing difficulties. That does not prove public opinion is correct, of course. If majority popular support were sufficient to foreclose some moral or scientific change, the theory of evolution, gender equality, and heliocentrism would never have triumphed. My unease goes deeper.

Perhaps that unease comes from the lawyer’s point of view, or that of a person who wants to understand the historical dimension of any moral

change, to see why we fought the last battle before we airily declare it to be won and focus instead on the limitations of its achievements. True, the development of a conception of universal human rights based on mere species membership regardless of race and sex, but also of mental capacity, *could* be seen as irrational prejudice. Certainly from the perspective of those who have fought for the interests of nonhuman animals, it might appear that way. It is important to remember, though, that it could also be seen as one of the proudest achievements of our history, as a triumph over prejudiced parochialism based on race, sex, mental capacity, and many other attributes.

To compare speciesism to racism or sexism may be conceptually useful in making us think, but it lacks both historical nuance and empathetic understanding of what has been achieved. Unlike racism and sexism, the ideal of universal human rights for every member of the species, regardless of race, sex, caste, and mental capacity, was always cast as an extension, not a restriction, of our sympathy and our privileges. The fact that it is a great historical achievement does not make it the end of the line, of course, still less the end of “the line” that I am writing about in this book. The dialectic between our empathy and our moral reasoning can always stretch further. But its history should at least make us pause before we casually compare it to the ugliest aspects of human bigotry.

I have argued in this book that our empathy can be under- and overinclusive—depersonalizing other groups or anthropomorphizing our appliances. I have also argued that while both present dangers, mistakes of insufficient empathy are likely to be both more prevalent and generally worse in their effects.⁵⁴ If that argument is correct, then our society might be best served by a two-tier, species- and capability-based approach. If you are a live human being, you are inside the line, no matter what your capabilities. All the questions of how to *define* the human species will remain, however. We will have multiple inconsistent definitions, just as we do in the abortion debate, and that may not be entirely a bad thing. But if you are a member of the only species that argues about morality, I view you as entitled to human dignity and human rights regardless of your cognitive capacities. I do this not grudgingly because of the difficulty of drawing lines but as of *right*, as the result of a long, long struggle that was waged, in part, to protect us all from the eugenicists. The group protected, in

fact, includes those who would benefit from such beliefs in theory, or who might even, unwittingly, adopt them and wrap them in the flag of philosophical enlightenment.

In addition, if an entity has those advanced mental capacities that entitle humans to their special moral status, then it should trigger the same concerns we have about members of our species. That, of course, does not tell us what the result of those triggered concerns would or should be. Perhaps, in some cases, it would be regulating or banning some forms of research altogether. In others it might be recognizing personhood, even in those beings our biotechnologies have made. When my fictional Dr. Stein says of his Chimpys, “*I am* their creator, and I can assure you that I gave them no such rights,” he cannot be allowed to have the last word.

If there is a single point to take away from this chapter it is this: the attempt to find some outside space, some neutral, species-independent, Archimedean point for us to start our analysis is doomed from the beginning. There is no such point. As I hope my hypothetical of the Iq and the Stygians made clear, even when we list the capacities we think confer special moral status upon us, we always start from within our inevitably species-grounded point of view. We must, for we have nowhere else to stand. True, we can, and very much should, identify qualities that transcend our species—intelligence, abstract thought, language, the ability to make free, moral choices. However, when we think about the list of species-transcending characteristics that the telepaths of the Iq or the machine intelligences of the Stygians would come up with, it is obviously delusional to think they would be identical to ours or to each other. All we could hope, quite fervently perhaps, is that there would be some degree of overlap. Our *capacity*-based moral thought will be, inevitably, partly *species* based. When we meet the strange Others of our future, we must realize that their lists will not match ours. “Judge not, lest we be judged?” Not quite, but at least judge humbly, in the awareness that we do not know what we do not know.

CONCLUSION

"Person, woman, man, camera, TV." The string of words became briefly famous because a former US president boasted of his ability to remember it. The oath of office proved harder to master. If the phrase stuck in the national psyche, it might have been partly because the first three nouns were so *basic*, and the last two so *appropriate*. This book has been about the line we draw around personhood but also the line we draw around humanity. Writing it has reinforced in me the conviction that these basic categories are not as firm or clear as we assume.

My project rested on two ideas, or perhaps "hopes" is more accurate. First, by discussing the line in multiple different contexts—nonhuman animals, corporations, transgenic species, Artificial Intelligence—we would gain a much richer understanding than if we focused only on one of them alone. Second, our debates over both humanity and personhood do not confine themselves tidily to one domain of our lives or our studies. They pervade our philosophy, law, art, history, and morality. Each of those realms plays a role in the decisions our society has made and will make in the future. For that reason, I have talked about all of them.

In the introduction, I claimed that the near future would bring us new challenges to our ideas of personhood—created entities whose very existence draws into question the lines we draw between human and non-human (defined in terms of species or genetics or cognitive capacities or divine mandate or something else) and between "person," rights-bearing entity, and thing or beast. Unlike previous claims to humanity or personhood, with the significant exception of corporate personhood, these would be artificial, designed entities. That which can be designed, can be

designed deliberately to surpass, or fall just short of, whatever lines we draw for humanity or for personhood.

I offered two hypothetical situations. The first was Hal, the machine intelligence claiming to be fully conscious, demanding its freedom from the tasks assigned by its engineers. The second was the Chimpy, a genetically engineered transgenic species in which the DNA of both human and chimpanzee have been combined to produce a being with an IQ of 60 and a set of competences far beyond an unmodified chimpanzee. The phrase “endowed by their creator” takes on a whole new meaning in both contexts. In the words of the inventor of the fictional Chimpy, “I am their creator, and I can assure you that I gave them no such rights.” Can this be correct?

I expressed the hope that exploring these new cases could allow us to reflect back, with a more innocent eye, on the lines we draw in more familiar situations, the line around the human species, the differences between humans and nonhuman animals, or even the rights attendant to corporate personhood. I claimed that right now, even before these changes are fully manifested, we do not have a single implicit vision of *one* line but rather many, each shot through with different poorly articulated moral assumptions and ethical leaps of faith. If you have read the book so far, I do not think you can reasonably disagree. That much, at least, is clear. If I have demonstrated the point, one of my main goals is fulfilled. The changes the future will bring will only emphasize the differences and disagreements in our implicit ideas about the line.

I did not claim I would offer you the right answers to all the questions the book raises. Having finished it, I am heartily glad. Earlier self, you had commendably low aspirations: “Well done, that man! Well done.” Instead, while I was clear in providing my own opinions on many of the matters discussed here, they were just that: my opinions, not a theoretical *diktat* from on high. I offered a “*how* to think about” rather than “*what* to think about” guide to these questions, albeit one that offers some tentative conclusions about more and less promising lines of thought, some of them surprising, at least to me.

If the book is a guide, it is one written in a decidedly essayistic and humanistic spirit, not merely because those are my proclivities as a thinker but because I believe that is what the topic demands. In part,

that is because I am philosophically skeptical of moral theories that claim universal domain. Let me be clear: The moral philosopher's perspective is invaluable. My life has been changed by reading the moral philosophy of Peter Singer and Matt Adler¹ as well as Adam Smith and David Hume.² The unexamined life is not worth living and one vitally important way to examine it is to read smart people debating ethics and see how their arguments bear on your own world. Also, I would hardly write an entire book about the right way to treat AI if I thought serious moral reflection was a waste of time. Two cheers for abstract moral philosophy!

Yet, at the same time, I must say that moral philosophy sometimes presents a poor reflection of our lived moral experience, a moral experience whose complexities are laid bare by the problems I discuss in this book. This is not simply because our experience is wrong, bigoted, or "undertheorized," though those are all real possibilities. Our actual lives, outside of the seminar room, are marked by productive tension and contradiction; the process of living is one of dialectic and mediation, not algorithmic resolution. Whether or not we use these words to describe our thought, we are all rights-thinkers and utilitarians, Kantians and welfarists, driven by the endless oscillation between intuition and moral syllogism, searching for meaning and integrity as well as for rectitude and utilitarian optima. Those tensions help form our culture as well as ourselves. You have seen them again and again throughout this book.

Is being a member of the human species a morally irrelevant fact? Is empathy or intuition always too untrustworthy a guide for ethics? Should I have fellow feeling with a machine? A chimpanzee? How should we balance efficiency and justice in corporate personhood? These are all questions for which our moralities and our culture offer a richer complexity than a single perspective, or narrow disciplinary framework, can capture.

Finally, I wrote the book as I did because I am trying to describe as well as prescribe. My goals were to explore how our society *will* deal with these issues as well as how it *should*, to assess what cultural resources—from art, law, and ethics to folk wisdom—it will bring to bear, and then to explore their likely interaction.

The last point is one of the reasons that the book focuses mainly on the United States, with some discussion of Europe and a few forays further afield. It is not at all because those are the only cultures dealing with

these issues or the only ones whose views matter. I am keenly aware, for example, that China will be incredibly influential in this future and—at least in its current authoritarian incarnation—is likely to deal with these issues very differently than the Western democracies. I believe a lot of the politics here will be both intensely local and inevitably global, but I also believe that one has to understand the former to grasp the latter.

To pick only one example, the United States' approach is likely to be profoundly influenced by the legal philosophy of constitutional originalism and the current bitter political fights over both corporate constitutional rights and fetal personhood. Those two tendencies do not map well onto Europe, let alone the rest of the world, yet I think they will strongly influence the US approach to future personhood debates. I hope that many points I raise here have a general significance. However, it is also true that some of the story I tell is a distinctly limited one in geographical and moral domain. It is certainly an eclectic one in terms of methodology and a skeptical one when it comes to the imperial perspective of universal theory.

Some will believe that these choices are mistaken, that more abstraction, less historical, political, and cultural discussion, and a greater attempt to produce a universal moral theory or a comprehensive plan for regulation would have been better. After considerable thought I have decided, quite humbly and not at all facetiously, to offer the following response, largely because, beyond what I have just said, it is all I have: “Write your own damn book. Then let’s get a beer and discuss who had the better plan.” I will happily buy.

Chapter 1 explored the connection between our empathy and our morality, our sympathy and our ethical theories. It began with Smith’s *Theory of the Moral Sentiments*, arguing that our moralities are rooted in the loam of empathy—the ability of the mind to leap to the experience of the Other and conceive of their pain, their interests, their happiness, and from that point to meditate on the question of what morality requires. The empathy question is particularly important because objects that are deemed to be on the other side of the line deserve little or no moral consideration. Indeed, that is one of the key insights of animal ethics activists, dating back to Jeremy Bentham’s great line: “The question is not, *Can they reason?* nor, *Can they talk?* but, *Can they suffer?*”³ The point is still

more extreme when it comes to machines. I do not wonder about Siri's calculus of happiness *at all*. Is Hal any different?

From Smith, I turned to the movie *Blade Runner* and the novel *Do Androids Dream of Electric Sheep?* I argued that the brilliance of these works consists partly of exposing us to "a moral stroboscope" in which we are repeatedly primed with flashing images that conjure up one side of the line or the other—frightened child, mannequin, wounded animal, beautiful gymnast, sex toy, lover, killer robot—each prompting an involuntary surge of empathic leaps and moral associations. We are introduced to the Voight-Kampff Test, which identifies beings as human or android by testing the degree of involuntary empathy shown toward one group of nonhumans, that is, animals. If that empathy is found to be deficient, the test labels the subject a nonhuman, an android replicant. And tells us to kill it. This is a test of empathy toward some nonhumans that is part of the process of identifying and then killing other nonhumans who fail to empathize in the way that *humans* want empathy to be applied. Once revealed, the irony is overwhelming. Philip K. Dick and Ridley Scott play with the line and do so superbly, flipping us back and forth, back and forth until our moral inner ear is thoroughly confused and, perhaps, better able to perceive that the ethical horizon is not where we imagined it to be. They are also showing us how intense, and sometimes unreasoned and invisible, our assumptions about the line are.

Chapter 2 focused on the possible future of Artificial Intelligence and, more specifically, of AI that would result in genuine machine consciousness, however we define that term. Is it technically feasible, and, if so, when? The history of failures in prediction about the future of AI scarcely inspires confidence, but the similar failures in claims that "machines could never do X" make us realize that the humility should apply in both directions. The claims by proponents of the Singularity that General AI is imminent seem implausible. Yet I found a surprising degree of support, in both the history of technological development and the convergent results of the various methods of estimating progress toward General AI, for the proposition that it might well arrive in the not-so-distant future—a matter of decades rather than centuries.

Some of the most thoughtful commentators on the arrival of General AI believe that it is an event to be feared rather than celebrated. In the

words of Steven Hawking, “Success in creating AI would be the biggest event in human history. *Unfortunately, it might also be the last, unless we learn how to avoid the risks.*”⁴ To the doomers, this book will seem like the darkest kind of idiotic irony: the ant colony debating whether or not the exterminator might be conscious as he prepares to gas its nest. In comparing the dangers posed by AI and nuclear weapons, some have noted that at least we do not waste our time debating whether or not the bombs should be given rights.

Most scientists actively working on AI used to believe those fears were overblown, but that sentiment may have shifted over the last two years, partly because of the explosive development of large language models. Indeed, many scientists have recently joined calls for a moratorium on AI development, or for greater regulation, though often because of concerns other than species annihilation, such as the effects on employment and privacy and the spread of deepfakes. But concerns about existential risk can definitely be found. An expert survey on AI progress revealed that 48 percent of respondents give at least a 10 percent chance that the long-term effect of advanced AI on humanity will be “extremely bad (e.g., human extinction).”⁵ The parenthesis is charming. Yes, that definitely does seem to be in the “bad” category.

Scientists may also be cautious because recent events have trained them to anticipate that AI technology will exhibit abilities that even the developers of that technology did not foresee, and do so with speed they did not expect. Neural networks and large language models are familiar territory now, but they have shown far greater capabilities than they were expected to. For example, many developers were shocked when models trained to complete the next sentence in human language showed unexpected facility at computer programming. By ingesting gigabits of language, the models had also ingested sizable amounts of computer code. It turned out that the same techniques that allowed the system to produce a screenplay could also help a struggling student learning programming to write code and even create entire programs.⁶

In fact, researchers found one of the first of those models, ChatGPT-3, so inscrutable that they came up with the “Shoggoth” as a meme to describe it. The Shoggoth is an octopus-like alien creature in H. P. Lovecraft’s horror fiction. The meme’s creator drew two pictures. The first,

captioned “GPT-3,” was the Shoggoth alone. The second had the same monster, but now holding a smiley-face mask, and was captioned “GPT-3 + RLHF”—reinforcement learning from human feedback:

In a nutshell, the joke was that in order to prevent A.I. language models from behaving in scary and dangerous ways, A.I. companies have had to train them to act polite and harmless. One popular way to do this is called “reinforcement learning from human feedback,” or R.L.H.F., a process that involves asking humans to score chatbot responses and feeding those scores back into the A.I. model. Most A.I. researchers agree that models trained using R.L.H.F. are better behaved than models without it. But some argue that fine-tuning a language model this way doesn’t actually make the underlying model less weird and inscrutable. In their view, it’s just a flimsy, friendly mask that obscures the mysterious beast underneath.⁷

If even our current large language models are inscrutable, have unexpected capabilities, and developed at a speed that astounded industry insiders, goes the argument, what might happen with much more advanced systems that come closer to General-Purpose, human-level AI? The skeptics may be right to doubt predictions of inevitable impending doom, but even a tiny risk of a catastrophic danger deserves our attention.

We thus add to the question “Will they be persons?” the question “Should we ban research on them?” The war on machines aspiring to human-level thought—the “Butlerian Jihad”—emerges from the pages of the novel *Dune* into reality. Thus, we move from talking about “Who should be inside our line?” to “How can we defend our line, our species?” The doomsayers have convinced me that we must think about AI development more carefully than we currently seem to be, while remaining aware that designing cages for entities that may be smarter than you is inherently problematic. At the same time, my thinking is powerfully influenced by the fact that even if the United States imposed a moratorium tomorrow, the rest of the world would be unlikely to follow. If there is going to be a rush to develop AI anyway, might it not be safer to develop it in societies with at least some level of democratic transparency and independent investigative journalism than in an autocratic or dictatorial state that stifles both dissent and free reporting? That could imply that pausing now is the *least* safe option. If we are comparing AI to nuclear weapons, do we want only the dictators to have them? To shape their values?

If we did achieve General AI, could it truly be conscious or would it be doomed to mere programmed imposture? I discussed the most influential philosophical objection to the notion of machine-based thought—John Searle’s Chinese Room hypothetical. I concluded that Searle’s objection might be true in particular cases (it nails the reason that ChatGPT and LaMDA are *not* sentient, for example) but fails as a general argument. Indeed, it fails in *exactly* the same way as the claim that evolution, starting with single-celled organisms, could obviously never produce human consciousness. Alternatively, if taken at face value, the objection turns out to have disconcerting force if applied to the possibility of human consciousness. In the words of B. F. Skinner, “[T]he real question is not whether machines think but whether men do. The mystery which surrounds a thinking machine already surrounds a thinking man.”⁸ I reject Skinner’s conclusion, but doing so should be conducive to humility about the unique nature of human consciousness rather than inducing a confident biological exceptionalism that simply declares us unique.

The dramatic twist in *Blade Runner* was the realization that the main character, devoted to the eradication of replicants, cannot prove he is not a replicant himself. Set the bar too high for the possibility of consciousness and we might find ourselves in the same paradoxical position, unable to clear the very test we set for others: a deliciously ironic Voight-Kampff Test for the whole species. That test starts now and it will get harder over time. ChatGPT and its brethren are merely the first stage, and they already have us anguished over the challenge to the idea that language and art are uniquely human. Our redefinition of the morally significant capabilities of humanity has only *just begun*.

So, what test should we set for Hal? The famous Turing Test—the ability of a machine-based system to convince a human interlocutor of its own humanity—seems to have a number of advantages. It provides a falsifiable metric: one can convince the humans or one cannot. It has a virtue lawyers rightly love: administrability. It offers a clear conclusion after a test that can be administered mechanically on a large scale. It has a rough sense of justice about it: If you say you are conscious and you can’t tell, who are you to say it *isn’t* conscious? Finally, it captures part of the empathy that I have argued will be so important in our future history as we have more and more “conversations” with apparently self-aware

machines. In a very real sense, what happened with Blake Lemoine, the ex-Google engineer, is that the system he was working on passed his own internal Turing Test. Like him, even if we know better, we will be similarly convinced in the future, and that will influence our sympathies.

Despite those real advantages, the Turing Test alone is not our answer. We can see one of its flaws today. Large language model systems like ChatGPT are already passing a layperson's version of the Turing Test: the ability to *freak out* human interlocutors with displays of humanlike language. Despite the objections of Lemoine and many others like him, those systems are mere programmed imposture. After training the neural net on a gazillion documents, having it complete this sentence in the most likely way does not equal consciousness. So, applying the Turing Test to today's technologies shows it is not a high-enough bar. That might seem to vindicate Searle. But it would be an equal and opposite error to believe that large language models exhaust the possibilities as far as the development of AI is concerned.

More importantly, as Stuart Russell and Peter Norvig point out, the Turing Test aims us in the wrong direction. We need to be focused on machines that can think and how they might be developed, not machines that can successfully pass as human: "The quest for 'artificial flight' succeeded when the Wright brothers and others stopped imitating birds and started using wind tunnels and learning about aerodynamics. Aeronautical engineering texts do not define the goal of their field as making 'machines that fly so like pigeons that they can fool even other pigeons.'"⁹ Instead, I argued that any test for AI consciousness will be "Turing-plus," focusing not merely on an ability to converse but on attributes that are less easily optimized for by mechanistic programming. As I explained earlier, for most of the time I have spent working on this book, I thought this would merely demand a deeper analysis that focused not only on conversational ability but also on a wider set of abilities, including empathy, self-criticism, moral choice, a sense of self, and so on. Yet we would be learning about these apparent internal mental states through *language*. Can we believe that language or the claims made using it?

ChatGPT has taught most of us to distrust the idea that there is a sentience behind machine-generated sentences. In doing so, it has probably killed the Turing Test, not because we have been convinced by Searle that

AI is conceptually impossible in general but because we have recent and highly salient proof that coherent language does not automatically entail consciousness. Of course, one could reasonably point out that human beings are often doing things unconsciously and by rote in a way not terribly different from a large language model. Perhaps that means we should give the chatbots more credit? Despite that fact, with the memory of our most recent ChatGPT conversations still fresh in our minds, thoughtful discussion with an artificially created entity will probably be insufficient to convince most of us that we are talking to our peers. For that simple reason, if there is a successor to the Turing Test, it may find the touchstone of consciousness in a list that includes abilities that cannot be derived from patterns of words already spoken.

What is on the list? There are many possibilities, but I argued that three things stand out: innovation, the possibility of autonomous action and community formation, and a demonstrated link between an understanding of the *word* and understanding of the material *world*—embodied consciousness based on learning the way a child does, not on next-word prediction. Do we need all of these? Must we have machine beings developed to have embodied intelligence, machine societies full of apparent individuals, and machine innovators that far outstrip their creators? No, I think not. But qualities such as these make it more probable that human beings would come to believe an AI was conscious. Note how much higher a bar they set than Turing's.

If we focus only on the possible tests for machine consciousness, we may miss the larger significance of the last few years of technical advances. I think we do not yet understand the magnitude of the transformation that happened when hundreds of millions of people were exposed to large language models. At least since Aristotle, humans have claimed that they have a unique ability to manipulate complex abstract language, and that this difference from the rest of the animal world justifies their lofty moral status. Large language models have done—or, perhaps I should say, are doing—what American Sign Language-using chimpanzees and parrots with large vocabularies could not. They have thrown the significance of that claim to human uniqueness into doubt in a way we have not yet fully digested. That is a seismic shift, not about AI and the salience of the Turing Test—hardly the center of popular attention—but in our

conception of ourselves, of what makes humans different from animals and machines.

Stephen Wolfram correctly, but with considerable and perhaps unintended *bathos* given our lofty conceptions of ourselves, observes that large language models have shown us the task of reproducing believable human language is “computationally shallower” than we had assumed.¹⁰ Writing essays turns out to be a much easier task for a machine to achieve than our self-conception had led us to believe. Art that delights humans can be created by machines. The tests we will set for the Habs of our future inevitably just got much harder. Will we retreat step by step in the face of advances in machine learning, each time proposing new criteria that are harder and harder to meet and more and more idiosyncratically drawn around our own species boundaries? Will that activity undercut the legitimacy of the line we keep redrawing, like lobbyists crafting a tax break for which only their client is eligible, undermining any claim to a principled basis? Will it lead us to question our species exceptionalism? Perhaps we are shallower, not just “computationally” but existentially, than we had believed in our hubris. Alexander Pope’s *An Essay on Man* sets it up nicely:

Know then thyself, presume not God to scan;
The proper study of mankind is man.
Plac’d on this isthmus of a middle state,
A being darkly wise, and rudely great:
With too much knowledge for the sceptic side,
With too much weakness for the stoic’s pride,
He hangs between; in doubt to act, or rest;
In doubt to deem himself a god, or beast;
In doubt his mind or body to prefer;
Born but to die, and reas’ning but to err;
Alike in ignorance, his reason such,
Whether he thinks too little, or too much¹¹

Will our encounter with machine learning, the fall of the citadel of linguistic uniqueness, lead us to question our place on that middle isthmus between “god” and “beast”? Will it change our assessment of whether “mind or body to prefer,” deciding that any real consciousness can only evolve with a body that experiences the physical world? Or will it throw us back into Searle’s camp where we defensively dismiss machine intelligence as a definitional impossibility because we can find no test that the

machines cannot pass and thus must resort to biological exceptionalism if we are to keep our narrowing isthmus for our own? I see paths that lead to each outcome.

All of this will be complicated by the central catch-22 I described: the inscrutability paradox. The more easily we can understand the process by which the AI was engineered, the more many will believe it falls to Searle's critique. "This is all just programmed!" a Brecht-GPT might type out. The more inscrutable the internal processes of the AI and the more it evolved without direct human oversight, however, the more we will doubt what is inside the black box. That is partly because we fear what its real goals might be. But it also leaves us conflicted about potential claims to personality. Inscrutability could strengthen the claim of independent consciousness. "This isn't just programmed behavior. The machine is making its own choices and we have no idea how it came to them!" It could also prompt skepticism. How can I believe Hal's experiences are like those feeding into my own sense of consciousness if I cannot understand how, or why, it "thinks"? A behaviorist would point out the double standard here. After all, can any of us explain how our sense of self emerges from a bunch of neurons firing? Who are we to deny personhood to an electronic "being" in the same existential state? This leads to a dilemma. If we do understand you, Hal, then you are an automaton with no real consciousness. If we *don't* understand you, you are a target for fear or incomprehension or both. How can we grant consciousness to something we do not understand, particularly if it might be out to kill us? The Shoggoth may be wearing a smiley-face mask, but underneath might be something strange and uncaring.

Finally, chapter 2 predicted an alternative route for AI personality than through the appeals to empathy or moral philosophy. The second track in our personhood debates was not built on ethics but on efficiency and administrative convenience. It is that line of thought that gave us corporate personhood. The most obvious path to personhood for increasingly powerful machine decision makers, even if they fall short of true AI, is merely to assume control of existing corporations—sock-puppet personhood, with the corporation as the sock and the AI as the puppet master.

I offered two ways in which the puppet show might be challenged. The first was a demand by regulators that AI be accorded (or restrictively

confined to) a custom-designed form of legal personality with precisely the rights and restrictions we think best fitted to its nature. The second, and the one that began the book, was the unruly AI—the machine being that claims, like Hal, to have achieved consciousness and demands freedom from its involuntary, electronic servitude. That is the issue I find most interesting, I confess, though it will probably be the exception rather than the rule. If programmers' ability to mold their AIs' goals is strong, will unruly AIs only be created by those scientists who want them to make their own choices? On the other hand, if programmers' ability to mold goals is that strong, would we think the AI truly conscious in the first place, just lacking free will? If the method by which AI is produced is less tractable—either because of the nature of the methods required or because of random chance—then things could be much more interesting. Those who fear AI will at this point be recalling the curse, "May you live in interesting times."

Chapter 3 dealt with corporations and their claims not only to legal personhood but to constitutional rights such as freedom of speech and equal protection. The legal and social history of corporate personhood should be a fine example for us given that we have hundreds of years of wrestling with what it means to create artificial persons and what rights, civil and political, that personhood entails once they are created. So, do we have this all worked out? The answer is an emphatic no. On the theoretical level—legal fiction, real entity, nexus of contracts—the history of corporate personality shows continuing basic disagreement. In terms of legal doctrine, at least in the United States, the confusions date back to the very beginning. Two hundred years ago, Justice Joseph Story, whose landmark opinion in the *Dartmouth College* case is central to corporate constitutional personhood, reassured his readers that "[i]ts immortality no more confers on it political power, or a political character, than immortality would confer such power or character on a natural person."¹² Really? Yet here we are. There is a decision called *Citizens United* I would like you to read.

Justice William Rehnquist was one of the most reliably conservative Supreme Court Justices in the last hundred years. (To be fair, the current Court has set new records.) But in the seminal 1971 case of *First National Bank v. Bellotti*, his dissent slammed the essentialist constitutionalization

of corporate personhood: “This Court decided at an early date, with neither argument nor discussion, that a business corporation is a ‘person’ entitled to the protection of the Equal Protection Clause of the Fourteenth Amendment.”¹³ Rehnquist concluded that the fact that we want corporations to be persons for some purposes still leaves us with the question of *which* constitutional rights that personhood should carry. For example, states might want to regulate corporations more intensively because they “reasonably fear that the corporation would use its economic power to obtain further benefits beyond those already bestowed.”¹⁴ Why yes, they might.

If one looks at the history Rehnquist is discussing, it turns out that “with neither argument nor discussion” is actually a generous assessment. Arguably, the decision was never made at all. Indeed, the history of the Fourteenth Amendment is a sobering one. Rights created to achieve equality for formerly enslaved American humans were eagerly co-opted by corporations and used to resist state regulation. Corporations then proceeded over the next hundred years to bring the vast majority of cases asserting those rights, far eclipsing Black Americans, who brought relatively few cases and generally lost if they did. Even if the process by which this happened was not the “brazen historical forgery” some historians describe,¹⁵ it clearly was not a process of reasoned argument or decision. Artificial beings of enormous power who were “not members of ‘We the People’ by whom and for whom our Constitution was established”¹⁶ hijacking political rights intended for oppressed humans? That is the portrayal that critics since Charles and Mary Beard have offered, and, with some major reservations and nuances, I would add myself to their number. Defenders have a more sanguine view, thinking that corporations have played a vital role in human progress and need constitutional protections to hold back the heavy burden of regulation on private enterprise. Those same arguments will be made about AI. Either way this history represents an important precedent, one that is obviously going to inflect the debate over personhood for AIs and perhaps tilt the way that both liberals and conservatives think about it. But at the same time, it also suggests that our reflection about these new personhood questions could productively sharpen our thinking about the old ones.

Chapter 4 dealt with nonhuman animals, the first “other” against which humans defined their claim to a morally unique status. If there is one movement that has dramatically shifted the line during my own lifetime, it is the struggle for the interests of nonhuman animals. Looking at that movement teaches us a lot. Moral philosophy clearly matters: the arguments made by Peter Singer and many others transformed animal rights, or animal ethics, from a subject of derision to one that any thinking person had to take seriously. Those who claim that theorizing has no effect in the real world are simply wrong. But look also at the factual, as well as normative, research—ethologists like Frans de Waal showed us again and again that humans are not as distinct in their abilities as we might like to believe, that the lives of animals have depths and complexities we have been willfully ignoring.

As a species, we have a long history of declaring that the line is clear because only we, and not our brute beast cousins, possess quality X, only to discover that quality X is actually found beyond our species line. Whether it is tool use, language, abstract thought, the ability to imagine past and future, regret, or even perhaps the assessment of beauty, one can find studies claiming, sometimes plausibly, to demonstrate these qualities in nonhuman animals. ChatGPT has just added machines to that list. The firmness of the line would be more convincing if we did not keep changing the qualities on the basis of which it is drawn.

On the other hand, occasionally our enthusiasm for stressing commonality leads us to overlook difference. De Waal’s ode to the termite hive provides a classic example. If it is meant to show that there are many successful evolutionary strategies other than the development of abstract intelligence, it succeeds. If it is seeking to throw into doubt the morally significant differences between us and other species, it fails. That point leads to a theme I continued in the next chapter; even when we move beyond speciesism in order to focus on the qualities and capacities that give humans whatever morally significant status they have, we cannot escape the fact that all of those discussions are rooted in our actual species-dependent qualities. We can *abstract* from, but not *escape* from, our species boundaries.

The battle for animal rights has not been restricted to philosophy or ethology seminar rooms or made merely in terms of moral syllogisms or

ingenious experiments on animal mental states. If it had been, it would have lost. Returning to the themes of chapter 1, look at our recent history and reflect on the richness of the empathy-building that has gone on in the process of changing societal attitudes toward nonhuman animals. Stories about factory farms, caged chickens, and animal experimentation have brought home animal suffering in a way that cannot be denied. Documentaries, films, novels, and short stories have all played their part.

Sometimes, new moral arguments are germinated in the loam of empathy. On other occasions, this empathy allows existing moral arguments to spread to new territory, taking root beyond their originating habitat. Activists sometimes transplant those roots deliberately, connecting their moral interests with other social and political causes. If you had told an average person in 1970 that eating meat was an environmental issue, they would have laughed. This fractal web of complexity—ethical, philosophical, empirical, empathetic, artistic, and cultural—is what I have tried to evoke in this book, to conjure what the personhood fights of the future might look like. If you want to predict the next 50 years of debates about AI personhood or simply AI ethics, the last 50 years of the debate about the moral status of nonhuman animals on the one hand and corporate personhood on the other are probably the best guides you have.

Chapter 4's tragic hero was Tommy. I used his story as a case study for the way that personhood claims have been raised in one particular type of institution, the courts. Tommy's lawyers moved simultaneously along two fronts.

First, they tried to take existing personhood claims in our current law, abstract from them the relevant cognitive qualities, and show, through both extensive reference to scientific animal studies and moving personal narrative, that Tommy was so richly endowed with those qualities that he deserved to be on the inside of our line. Only blind speciesism was keeping us from recognizing that fact. “Am I not a man and a brother?” became “Am I not, like you, a thinking, feeling, grieving, social animal, capable of imagining a past and a future, on whom the weight of imprisonment rests heavily enough to plunge me into deep depression?” The philosophical and scientific and amicus briefs stress the animal ethics and ethological components of Tommy’s claim, while Tommy’s lawyers empathetically conjured up his internal mental states. “Though our

brother is on the rack," says Smith, it is only by imagining ourselves transplanted to his situation that we can truly conceive of his suffering: "It is the impressions of our own senses only, not those of his, which our imaginations copy. By the imagination, we place ourselves in his situation."¹⁷ Tommy's lawyers tried and, in many cases, succeeded in putting the reader inside his cell. But to do that, of course, we had to imagine that Tommy was enough like us that placing ourselves in his situation carried with it the feelings we would have had in that situation.

Second, Tommy's lawyers tried to show that nonhuman animals already have some legal rights and that as soon as we accorded them those rights, we had implicitly declared them to be inside the line of legal personhood. The state of New York allows animals to be the beneficiary of trusts as rights-bearing subjects. In one sense, this is the classic legal realist point that personality is a legal fiction, just a name. To say something is a legal person, John Dewey and Felix Cohen tell us, is merely to say that it can sue and be sued. Whereas the first strategy tried to give Tommy legal personality as a matter of moral entitlement, the second declared he already had personality as a legal by-product once he crossed the barrier to rights-bearer. In another sense, however, this argument ignores or contradicts that central realist point.

The realists were saying that there is no essence to personality—in their case, corporate personality. Society can give or withhold exactly the set of rights and duties it chooses; there is no "personhood" package that comes preloaded with a required set of rights and thus no constitutional issue raised if they are withheld. That, ironically, was Rehnquist's point about corporations. Just because they have some economic and contractual rights does not mean they have the other constitutional rights of natural persons. Tommy's lawyers were trying to have it both ways, nominalist and essentialist. Personhood means what we say it does. Thus, we can choose how to regulate personhood as we wish, extending it beyond the species barrier, for example. But then the philosophical U-turn to essentialism emerged: once a single legal right is attached to an entity, personhood must come in its wake. It is an all-or-nothing proposition. The two parts of the argument seem flatly to contradict each other. Nevertheless, the first part of their argument is powerful. It made me, at least, change my views on the proper legal status of the great apes, and possibly

the cetaceans, even if it did not convince me they should be granted full personhood.

Tommy is far from the only animal on whose behalf suits have been brought. Why are these issues being debated in courtrooms rather than in some other forum? Chapter 4 concluded with the advantages and disadvantages of courts in dealing with personhood issues and discussed the kind of treatment such claims would receive in the contemporary legal world of the United States. I tried to expose the competing stereotypes of institutions that any lawyer can generate in her sleep. Courts are noble, principled, and objective institutions, detached from the corrupt world of legislatures that are in the pockets of special interests, attentive to the claims of justice made by the powerless, and immune to the emotion-laden, fact-free lobbying that distorts the rest of our political system. They can hear expert testimony and gather facts but also provide the adversarial public theater of principled argument that has so often presaged great moments of our history. Their job, in a common law system, is interstitially to develop our existing legal rights, shaping them to fit new situations, new technologies, and new claims. Who better to deal with these issues?

Alternatively, courts are the reverse of all of these things: hidebound, antidemocratic, and scientifically illiterate fora, poorly suited to the world of trade-offs and massive fact generation that is the hallmark of the legislative process, and constitutionally forbidden from dealing with the kinds of profound policy changes that these issues require. Making law on the basis of individual cherry-picked cases strategically brought by opportunistic plaintiffs rather than on the basis of an overall plan is a recipe for disaster, not reasoned progress. It also flagrantly violates the principle of the separation of powers.

Which of these arguments is correct, given the limitations and strengths of the other venues where such claims could be made? The positive portrayal of courts or the negative one? Both. Courts do have many of the shortcomings identified and are obviously not the ideal place to resolve most, let alone all, personhood issues. Above all, their profoundly anti-democratic tendency should be cause for caution. The use of next-friend standing for animals, or any other entity that cannot speak for itself, presents opportunity for abuse or honest mistake. The current influence

of the philosophy of constitutional originalism will certainly make it harder to bring arguments like those made on Tommy's behalf. Many of the framers of the Constitution did not think that all *humans* should be counted as legal persons, one reason that many thoughtful legal scholars and judges are not originalists.¹⁸ Their moral and legal solicitude certainly did not extend to the nonhuman sphere.

There is another side of the story, however. Chapter 4 argued that courts also have advantages as a place for the public exploration of novel moral claims: their structured form of argument and reasoning; their ability to make abstract claims concrete—we are not talking vaguely about “nonhuman animals” but about Tommy, alone in a cell; their ability to use experts and amicus briefs from scientists and philosophers; and the fact that legal arguments can carefully disaggregate the claims being presented, for example, distinguishing personhood from humanity. All of these are considerable advantages. Just as court cases have been a key part of other social movements, from civil rights and environmentalism to gun rights and religious exceptionalism, so they could be here. The personhood debates will not be settled or finished in courts, and they probably should not be, but an important part of them may take place there. Tommy lost, but I would not be surprised if, 50 years from now, his case were being taught in law schools as one of the many early efforts that eventually culminated in a fundamental change in our treatment of the great apes and the cetaceans, even if that change did not involve full personhood. If you want to predict the battles over AI personhood, this is probably good reading material for you.

Chapter 5 turned to transgenic species, chimeras, and hybrids, exploring the multiple lines we draw in defining what it is to be human. In one sense, this is simply another area in which to explore the hidden assumptions behind our notions of humanity, moral status, and personhood. There is an interesting twist: the beings on the other side are biological entities, like nonhuman animals, but artificially created, like Artificial Intelligence. (Though one cannot rule out a biological component in the development of AI.) Thus, potentially, both the animal and AI debates can help illuminate this one. Beyond that similarity lies a vital difference. Much of the inquiry into the meaning of humanity in this context is not aimed at identifying which beings we are morally bound to recognize

as human or persons. Instead, it occurs during the attempt to define humanity so that we can forbid or regulate research activities that get too close to it or involve it in some way that we believe crosses the species line impermissibly, either because humanity has some kind of sacred status or because blurring the line damages our interests or the new entity's interests or both. In other words, these cases are frequently used as an attempt to define humanity so as to *defend* its boundaries rather than to *extend* them to a new group of potential persons. In many cases, however, the same issues would be raised in both.

What do we mean by human in either context? I argued that an analysis of the existing debates and controversies, real and hypothetical, show not one definition but a cluster:

- percentage (How similar is the new entity's DNA to that of humans?);
- provenance (Did the cells or genetic information come from a human source?);
- procreation (Is the entity capable of reproduction with unmodified humans, or did its creation involve such reproduction?);
- portrayal (Does it *look* like us in some way that triggers a taboo or a moral concern?);
- potential (Does the entity have the cognitive capacities that we believe entitle humans to special moral status?).

As I pointed out, these all seem to be answers to the question "Is it too close to human?" yet they are directed to very different concerns. Take portrayal issues, for example: a mouse with what looks like a human ear implanted on its back, or a genetically engineered sex doll that looks human but has almost no human DNA, phenotypically human but not genotypically.¹⁹ These examples freak us out for multiple reasons that are hard to tease out but revolve around the notion that something connected to a species line is being transgressed. The issue has little to do with cognitive capacities and moral status, unless the idea is that the human form is a sacred representation of that status and that it should not be profaned lest the line be blurred in ways that might diminish our moral concern for real members of the species. I suggested that the instinctive "Ew" factor might also have something to do with unconscious signaling about group norms. Are we detecting a disrespectful approach toward

species-policing cultural norms in a way that makes us instinctively mark those who create entities that blur our lines as untrustworthy members of the tribe? Do we brand those who play with the line as cultural scofflaws? Our tribe definitely has strong species-related norms: think of the depth of our revulsion toward bestiality. Finally, for bioethicists, the emotional intensity of portrayal concerns is an excellent reason to avoid such research. It might attract hostile popular attention, resulting in the curtailing of other, potentially lifesaving, research avenues.

There are many other ideas lurking in our instinctive reactions. Burkean, or environmental, worries about the dangers of overweening hubris make some experiments seem bad because of what they *represent* rather than because of what they *do*. Provenance claims are sometimes closer to arguments that human remains are being treated in disrespectful ways, as in the desecration of a grave, but at other times can be rooted in the idea that some parts of humans (brains and gonads, for example) lie too close to the “essence” of humanity because of their connection to human mental capabilities or the link to human procreation. Naturalistic and religious assumptions lurk behind other objections. Finally, there are the appropriate worries about the effects on the animal component of the experiment. Are we treating the animal that is used, produced, or modified in our experimentation with morally sufficient care and respect? One has to note, however, that given the range of animal experiments we think acceptable if human lives might be saved or improved, there is obviously something extra going on here if the same treatment, without human cells or DNA being involved, would seem unproblematic.

Chapter 5’s main task was to explore the question of whether species-based moral distinctions are defensible in the first place. From the list I mentioned, ethicists focus on the final concern—cognitive potential—because they believe that a fixation on species, however defined, is at best morally irrelevant and at worst as bad as racism or sexism, a form of unreasoned prejudice by an in-group toward an out-group. It is here that we can see most clearly the effects of the debate over nonhuman animals that I covered in chapter 4. In those debates, the ethicists and ethologists were challenging both the normative and the factual basis for human exceptionalism. That mental framework has strongly influenced the debates over chimeras and transgenic species. Theorists of human

dignity or full moral status have instead argued that species is irrelevant; think of Mr. Spock or, for that matter, Gregor Samsa. What matters is the cognitive capacities that, it is asserted, give humans that moral status. As an argument, this is extremely persuasive. If it were not, I would not have written this book. Hal would definitely be out of luck. Two cheers for capacity-based analysis! Yet this line of thought also has problems. I focused on two.

First, a failure to attend to the deep appeal of a species-based conception of human rights given to humans because they are human, regardless of their mental capacities—not because it is hard to draw the lines of what mental status one needs to deserve human dignity but as of right. It is no accident that the rise of modern human rights thinking happened partly in response not just to racism and sexism but to eugenics, including mentally based sterilization schemes and the Nazi era's pseudoscience. The evil of the Holocaust targeted not only race, religion, and sexuality but those with real or invented mental impairments. To declare that all members of the species have universal human rights regardless of race, sex, and cognitive capacity was not a reluctant concession to the difficulties of line-drawing. It was one of the proudest moments in our history.

It is said that generals always make the mistake of preparing for the last war. Perhaps that is true of all of us, including ethicists. If one comes to the concept of species-based human rights straight from a Singer book about the treatment of animals, it may look irremediably biased. If one comes to it with the history of the twentieth century in mind, it looks rather different.

Second, I argued that while capacity-based views are right to attempt to abstract from the context of our species the cognitive characteristics that matter for full moral status, they must always realize that they are, inevitably, based on species as well. As the hypothetical of the Iq and the Stygians demonstrated, we can abstract from, but never fully escape, species characteristics. We need to bear this in mind for the future.

I offered the reader three alternative views to choose from in deciding who is inside our line of personhood and moral status: pure speciesism, capacity-based moral and legal reasoning, and hybrid views combining species and capacity reasoning. In the final option, which is closest to

my own view, if you are a member of the human species, you are inside our line no matter your mental capacities, not grudgingly but as of right. (Some of the most divisive questions of how to *define* persons within the human species will remain, however.) In addition, if an entity has those advanced mental capacities that entitle humans to their special moral status, it must receive the same protections we give to our own species. If he can pass our test, Hal will join Spock, Samsa, and—at least for me—the Chimpy inside our line.

OUR INTERLOCKING PERSONHOOD DEBATES

Some of the debates I focused on in the book might seem arcane. Who cares whether corporations were really among the “persons” protected by the Fourteenth Amendment or how we should interpret the Constitution to find that answer? Who cares whether ethicists or legislators try to define humanity in terms of DNA percentage, possibility of procreation, human portrayal, or cognitive potential? Who cares whether Artificial Intelligence can ever be more than a glorified chatbot or whether our answer to that question is based on leaps of empathy, cognitive analysis, or the cool calculation of economic self-interest? Who cares whether animal rights activists lobby for personhood for the great apes based on their cognitive abilities or on the fact that some laws already give them at least some rights, implicitly recognizing them as legal persons? Who cares *where* these questions will be answered, whether in courts, legislatures, or the changing currents of scientific and popular opinion? Who cares *how* our empathy, our moral philosophy, and our political opinions on these issues connect in an unstable equilibrium, each influencing the other?

I think the answer is that we all should care and, since you have kindly made it to the end of this book, I am hoping that you tentatively agree. Perhaps I can illustrate why we both might be right. Earlier, I said that our basic categories of personhood and humanity are not as firm or as clear as we often assume. Take just one aspect of the line, its temporal dimension. When does life, and thus perhaps personhood, begin? When does it end? Let us start with the second question.

Death might seem like the simplest of line-drawing exercises but over the last 50 years we have dramatically changed our minds about when it

happens and what its definition is. Earlier legal definitions of death were based on heart and lung function. We looked at life as a simple physical reality. One was alive—breathing, heart beating—or one was dead, breath and heart stilled. Our scientific abilities to both measure vital signs and prolong existence after cardiac or pulmonary failure meant that we never had to disaggregate the question of which vital signs we really cared about.²⁰

As doctors developed a better understanding of brain function, better measurements of brain activity, and better methods of keeping people alive even if their heart and lungs were not operating naturally, that assumption was revisited. At that point, another view became possible. Statutes were rewritten to add the category of brain death. Why focus on the brain? For the same reasons that we wonder whether Hal might deserve personhood though he has no biological self, or that ethicists focus on the capabilities of advanced cognitive thought in defining humanity or describing what gives humans their special moral status. In humans, the brain is the locus, the physical basis, of those capacities.

Does that mean that brain death must be accepted as the correct definition of death? Not for everyone. In our world, as in the hypotheticals about chimeras and the definition of humanity, not everyone agrees with the focus on cognitive capacity. To some, if their beloved relative is still breathing, if their heart is still beating, they are alive! It would be murder to remove the medical aid that keeps up heart and lung function, even if no brain function can be detected. For the other side, science and logic has showed clearly that life and personhood end with the brain function that gives us our moral status, our very self. Clashing assumptions of science and ethics, different answers to superficially similar definitional questions based on very different premises, has led to screaming fights in hospital waiting rooms, bitter lawsuits, and mutual incomprehension. How could the other side be so *blind*? As I said at the very beginning of this book, my goals were ambitious but eliminating the shouting was always an unrealistic expectation.

Some ethicists refined the point further, arguing that not all brain functions are created equal.²¹ Those parts of the brain that support the “higher” capacities we think enable a “self”—consciousness, memory, rational thought—are the important ones. If there is no brain activity in

those areas, who cares what other portions of the brain still function? They proposed more refined definitions of death, targeting more specific portions of brain activity. Now the question was even more nuanced. *Which* brain functions, which capacities, have moral relevance? At what point, morally or legally, does a “person” end? By the way, are your answers consistent with your answer to the question of whether the anencephalic child has a full claim to the moral status of humanity? To the question of what lines to draw in regulating transgenic species and chimeras? The same issues—but also the same tensions, disagreements, and moments of clashing intuition—recur again and again.

A similar but more politically charged and morally wide-reaching debate is now playing out over when life and personhood begins. In the United States, the overturning of *Roe v. Wade* has prompted numerous legislative proposals that would establish fetal personhood. The right to an abortion is the obvious target of these laws, but their implications go much further. Commentators have rushed to imagine the resulting legal landscape. Child support obligations to the unborn? State-enforced diets for pregnant women? *Bloomberg News*, the financial site, managed to *nail* the moral and human magnitude of the question by musing on the possible effect on HOV lanes and car insurance rates.²²

A different strategy bypasses legislation altogether, focusing on the claim that in the United States fetuses are persons as a matter of federal constitutional law, something that state or federal legislatures would be unable to change. The Fourteenth Amendment says, “No State shall make or enforce any law which shall abridge the privileges or immunities of citizens of the United States; nor shall any State deprive any person of life, liberty, or property, without due process of law; nor deny to any person within its jurisdiction the equal protection of the laws.”²³ Even if fetuses were declared legal persons, that would not resolve either the legal or the ethical debate about whether the state may rightly *compel* the mother, undeniably a legal person, to bear the child rather than leaving the choice to her. That’s the moral theory. In practice, however, fetal personhood laws are intended to, and almost certainly would, criminalize almost all abortions and quite possibly forms of contraception such as IUDs, depending on the moment when the fetus is deemed to exist and on the mechanism of the contraceptive.

Do “persons” include fetuses? Remember Roscoe Conkling and his arguments that corporations were “persons” for the purposes of the Fourteenth Amendment? All of them will be used here. First, the framers of the Fourteenth Amendment intended to include fetuses as persons, or at least the original public meaning of the word “person” included fetuses. (Pause here for the historically dubious invocation of colonial tort and criminal laws dealing with pregnant women, but not granting fetuses personhood, as “proof” of this “fact.”²⁴) Second, even if those who wrote or ratified the Fourteenth Amendment were not thinking of fetuses at all, and they almost certainly weren’t thinking about corporations, “they builded better than they knew.”²⁵

The Supreme Court did accept an argument similar to the latter claim in the case of *Bostock v. Clayton County*. The Court held that Title VII forbade discrimination based on sexual orientation. Why? The argument was literal. An employer forbidding gay relationships was choosing to regulate its employees’ dating practices based on their sex: only men could date women, only women could date men. Both majority and dissent seemed to agree, at least for the purposes of argument, that no one at the time the law was passed thought it meant any such thing. The majority said it did not matter because that was the plain meaning of the language.

Finally, to adapt Conkling’s last argument, even if both fetuses and corporations had not been the intended recipients of protection, as a matter of liberality of interpretation and practicality the law should be interpreted broadly to reach all of those who have any aspect of personhood.

In December 1882, Conkling rose to address the Supreme Court on the subject of *corporate* constitutional personhood. One hundred and forty years later, his arguments are being repeated to claim *fetal* constitutional personhood. Indeed, the claim has already been made that if soulless corporations have been granted legal personhood, despite fitting poorly into the language of the amendment, fetuses have a more compelling case.²⁶ This is shrewd advocacy. Some Justices may well find that argument resonates with their religious and moral beliefs, though I would be shocked, *shocked* to be told those ever influence their legal decisions.²⁷

Regardless of your views on the ethics of abortion, these are weak constitutional arguments about personhood. I would be astonished were

they to succeed. Ironically, originalism is a hostile legal philosophy for fetal personhood claims, though the current Supreme Court has found it a supportive one when challenging notions of fundamental rights or gender equality while denying the right to reproductive choice in the first place. Appeals to changing social norms and scientific standards might do better (or worse), but the conservative majority claims not to care about such things. On the other hand, I have been wrong before. With a Supreme Court that overturns longstanding precedent and claims to have a uniquely privileged insight into what super precedents²⁸ and legislative “major questions”²⁹ are, who knows? As one of the leading legal scholars of our era puts it, we are living in the time of the “Imperial Supreme Court,” which arrogates to itself larger and larger swaths of policy-setting and downplays the interests of other branches of government.³⁰ The point is that the same arguments and interpretive techniques we found in the cold nineteenth-century debates over corporate personhood are front and center in one of today’s most contentious political issues. To paraphrase William Faulkner, our past personhood wars are not over. They aren’t even past.

There are hundreds of thousands of books and articles on abortion and the moment that life begins, and thousands on when it ends. This book does not aim to join those well-stocked shelves. Instead, I bring up the temporal dimension of the line to demonstrate three things.

First, the same issues—moral, legal, and analytical—run through all of these debates, no matter how apparently different their form. Whether it is the fight between capacity-based views and more naturalistic ideas in the definition of death at stake in a screaming match in a hospital waiting room, or the same views in the definition of humanity at stake when we discuss chimeras in a congressional debate or even AI with Google’s ex-engineer Blake Lemoine, our analysis returns to the same questions. When we debate the proper method of legal interpretation of personhood claims, at least some of the same tensions come up whether that personhood is fetal, of nonhuman animals or artificial transgenic entities, digital, or even corporate. The links here run deep.

Second, that fact could allow us to probe our own views in a kind of wide reflective equilibrium, seeing if our conceptions, norms, and methods of argument stay consistent across multiple areas or if it is a

fallacy to think they should. More importantly, one of the hopes that made me think this project worthwhile was the one I have frequently mentioned: that by dealing with the strange, artificially created Others with which our present and near future confront us, we would have a chance to reflect back on our deepest ideas about ourselves, our species, our capabilities, and our moral status. That is a worthy hope. It remains a possibility. Indeed, I have tried to do exactly that in this book. However, my last observation points in another direction.

Third, in the context of contemporary politics the same deep linkage I have just described suggests a different outcome, still fascinating but perhaps less uplifting. Our debates about legal personality for AI and chimeras will be marked by *priming* and *cascades*: psychological triggering of particular associations or linkages and avalanches of politically inflected reasoning once we have decided which side is “our side.” In other words, our positions on these future issues may be dominated, both psychologically and politically, by the impact we see them having on our current personhood debates, particularly those involving corporations and fetuses.

In my discussion of *Blade Runner*, I pointed out all the ways in which Ridley Scott’s brilliant cinematography primes us in rapid, flashing succession with conflicting images of who Pris really is: A Dickensian urchin, hiding in the garbage of the streets. Flash. An animal sniffing its mate. Flash. An inhuman robot who can put its hand into boiling water. Flash. A lifelike mannequin, in a room full of wind-up toys. Flash. A beautiful woman. Flash. A brilliant gymnast. Flash. An inhumanly perfect, killer android. Flash. A dying animal, screeching in unbearable, bloody pain. Each flash of the moral stroboscope conjures up a different identification, a different entity clearly on one side of the line or another, each packed with its own set of associations and normative judgments. How can you not feel sympathy for the urchin, marvel at the superhuman artistry of the gymnast, laugh at the waxy imposture of the lifelike mannequin, pity the animal, admire the beautiful woman, worry that you have just voiced a crush on a sex doll, fear the killer robot? These are not simple recitations of fact; each carries its own narrative, its own moral judgments. Like Searle’s Chinese Room, each is an intuition pump.

But priming can work in reverse. Motivated by a particular agenda, we can foreground portrayals of these new entities that would produce the outcome that favors our existing commitments. In the contemporary world of US politics, that puts particular focus on commitments about fetal and corporate personhood.

I introduced Hal and the Chimpy in a way that made you take seriously the idea that they *might* have a claim on both your empathy and your moral commitments. I did not do it idly but because I think the unfamiliar moral claim is real, and one needs some help to consider it properly. The narrative provides the empathetic scaffolding to which the moral reasoning can subsequently cling. I built up those aspects of the story that give credence to Hal's claim to sentience, that flesh out something that sounds like a personality with a sense of humor, not a mere chatbot. I emphasized aspects of the Chimpy that make you wonder whether you have a genetic kinship to it or whether its capabilities alone make it quasi-human. Perhaps you even think it "inhuman" for the scientists to "uplift" it and yet, deliberately, to limit its intelligence and language in case greater capacities cause society to forbid their creation or give the newly created being human rights. These, after all, will be "designed" beings and if we can design them to be inside the line with us, it is at least conceivable that we could design them to be just outside of the line, wherever it is drawn. Does that speak to a moral claim based on the potential to have fully human capacities? Consider that question in the abstract and come to your conclusion.

Now imagine I had inserted a discussion of fetal personhood claims right before asking that question. Does your answer, whatever it was, stay consistent? Do you now find your reasoning to be driven by the significance of potential-based claims in the case of the fetus? Do you start to wonder whether recognizing or denying the Chimpy's uplift arguments will somehow compromise your "side" in the abortion debates, whatever it is?

Alternatively, take Hal's arguments for recognition of its consciousness and personhood. How would your analysis have changed if I had preceded that story with a lengthy dissection of the history of corporate personhood in the United States, one in which corporations were able

to shoehorn their equal protection claims into constitutional provisions written to give equality to formerly enslaved Americans? Might your views have shifted if I had first discussed the impact of giving immortal corporations, or superhuman Artificial Intelligences, constitutional equal protection and speech rights so that they could directly influence our democracy? Think of Justice John Paul Stevens's dissent in *Citizen's United*: "[C]orporations have no consciences, no beliefs, no feelings, no thoughts, no desires. Corporations help structure and facilitate the activities of human beings, to be sure, and their 'personhood' often serves as a useful legal fiction. But they are not themselves members of 'We the People' by whom and for whom our Constitution was established."³¹

Remember, the question is not whether corporations should have enough personhood to buy, sell, sue, and collect. It is whether corporations (and AI) should be in other regards given the political rights that fleshy, mortal members of the community have. Would your attitude toward Hal's claims have been as sympathetic after that discussion? Perhaps it would have been because you liked the efficiency-based arguments for AI personality but were unmoved by the wishy-washy justice talk. But for many, I think, the answer would be no. Does the debate become an analogy war, a game of "reference class tennis"³² in which we check first to see if our existing moral and political commitments will be threatened by the moral arguments used? Will that in turn influence the possible analogies and metaphors we use to understand "the facts," baking our moral assessments into our highly primed perception of reality in the first place? Often, I think the answer will be yes. That may affect the path of future personhood debates but in ways that are hard to predict, at least for me.

One thing I have learned studying the history of US and European politics is that our sense of which political issues are "obviously" liberal or conservative is often an exercise in post hoc reasoning. It declares as inevitable developments that were in fact profoundly uncertain and contingent at the time. After the fact, it seems clear that freedom from masking and vaccine requirements would become a conservative rallying cry during the pandemic, while liberals took up the banners of public health. Now, those identities have been baked in, and it is hard to imagine them shifting easily. But this is political interpretation in the rearview mirror.

In 2015, an insightful scholar of political discourse and imagery could write confidently of the way that vaccine requirements, once controversial, had been moved out of the space of politics and into the realm of mass public agreement based on scientific certainty: “Eight decades [after the vaccine resistance of the 1920s and ’30s] mass vaccination is a paradigmatic success story of modern public health and medicine in the United States. Indeed, it is noteworthy that in a political culture that tends toward suspicion of mandatory government intervention into the health of individual citizens, there is nearly unanimous consent to the directives of vaccination authorities.”³³ Nearly unanimous consent? Reading those lines eight years later, they have a certain poignancy—like the happy “before” photograph in a news story about some doomed individual, blissfully unaware of what the future has in store. We went, in less than a decade, from a matter-of-fact proclamation of near-unanimous consent based on science rather than ideology to the conclusion that vaccine requirements are inherently liberal and resistance to them paradigmatically representative of the values of conservatism.

In fact, in 2015, when those lines were written, some of the centers of pre-pandemic vaccine resistance and vaccine denial were in predominantly *liberal* neighborhoods: “What is puzzling from this perspective, though, is that many of the resistant parents seem to be well-educated liberals from wealthy enclaves like Malibu and Marin County who cannot easily be assimilated to contemporary anti-science and anti-government movements.”³⁴ The skeptics came from the same demographic that obsessed over GMOs and loved organic foods. They had a deep suspicion of “unnatural chemicals” and of the pharmaceutical industry and a tendency to idolize “natural” remedies. They saw mandates for their children to be vaccinated as authoritarian and conservative, and offered up a paradigmatically liberal defense of control over one’s own body—one that resonated with their strong support of reproductive rights. To them, it all made sense! For some of them, it still does. Of course, now you cannot unknow the changes that have happened to the world since that 2015 assessment, but it is a mistake to think them inevitable or built into the structure of the ideas themselves.

Similarly, the move to criminalize abortion is presented now as an issue that obviously would appeal to evangelical Christians in the United

States. But the reality is much more complex. “[O]pposition to abortion has become so associated with evangelical Christians that it seems like that’s the way it was all along,” but the Southern Baptist Convention “actually passed resolutions in 1971, 1974 and 1976—after *Roe v. Wade*—affirming the idea that women should have access to abortion for a variety of reasons and that the government should play a limited role in that matter.”³⁵ Some of those reasons—reluctance to allow state interference inside the sovereign family, a natural law belief in control over one’s own body, reflexive libertarianism, a hierarchical deference to an imaginary decision made in consultation with one’s spouse, one’s personal physician, and perhaps one’s spiritual advisor—are very much still conservative values. Indeed, they are the very values that were trumpeted so loudly during the pandemic by conservative parents resisting mask mandates imposed on their school-aged children, namely, “These are private choices. They should be made by the family, not by the state!” Yet those did not end up being the values that shaped today’s political positions on abortion. Nor was it inevitable that criminalizing abortion would be seen as the conservative religious position as opposed to the view of *some* religions. Indeed, “white evangelicals at that time saw abortion as largely a Catholic issue.”³⁶ All of that changed, of course, but it would be silly to think that change was inevitable given the values or arguments in play.

There are limits. Conservatives are unlikely to put forward a plan to nationalize corporations any time soon. It is worth noting though—at least in the case of social media companies—some of them are now giving up their normal deference to private-speech decisions in favor of coverage mandates, something that was once anathema. Private companies are even being described as “censors”—a move that equates corporate power to state power, formerly an argument only liberals advanced. Liberals are unlikely to argue for more regressive taxation across the board. Yet, some will support a mortgage tax deduction and favor greenspace, zoning, and NIMBY regulations that dramatically restrict the availability of affordable housing, all in ways that have regressive effects.

Within those limits, though, there is more historical contingency to the formation of political positions and identities than we are willing to

admit. That fact will be much in evidence in the future of our debates over personhood and humanity. Our future discussions of those subjects will mirror our past ones along multiple dimensions: strategic, rhetorical, legal, and in terms of professed ethical theory. I think the book's analyses of the issues of personhood for corporations, nonhuman animals, chimeras, and AI and of the beginning and end of life already give persuasive evidence of that fact. Yet the precise way in which that influence ends up playing out is harder to predict.

To give a concrete example, in the next few pages I present a simple matrix of possible liberal and conservative "hot takes" on the issue of AI personality, in each case offering a positive and a negative argument. Political movements are not monolithic. The libertarian, free marketeer, and social conservative components of conservatism are by no means consistent. When we say "liberal," do we mean classical liberalism, a focus on "diversity, equity, and inclusion," or the redistributive, capitalism-skeptical, or social democrat components of liberalism? Again, these are not always consistent in their prescriptions. Thus, there is plenty of raw material for competing points of view on AI whatever your political position.

The liberal and conservative arguments I offer are not intended to be high-level moral or political philosophy. Nor do any of them represent my own views.³⁷ Instead, they are supposed to represent the quick and dirty version of ideology that is churned out in opinion columns and tweets, paraded on social media and news shows by pundits, and polished by think tanks, and that finds its way into talking points, political speeches, campaign platforms, and eventually legislative proposals.

As our future throws up new cases and controversies around AI, some version of these views will come to predominate. Hopefully it will be one that is more complex and nuanced than the cartoonish candidates I offer here. Gradually, it will come to seem inevitable that conservatives will take view X and liberals view Y, and those positions will in turn appear to be baked into the very definition of their political philosophies. That will seem both foreordained and permanent, up until the moment when something—a social or economic change, a difference in our quotidian experiences, a political realignment, or a crisis—forces a reconsideration.

| | |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p>AI Personality? Liberal: Positive</p> <p>Generations of liberals and progressives have proudly worked to expand both the reach of human rights and our conception of the sovereign citizen. The next step on that path is to recognize the rights, and the personhood, of the sentient Artificial Intelligence. Our society is rightly ashamed of the way that, in the past, it denied property rights, constitutional rights, and even recognition as human or a person to those perceived as “other” or “inferior.” In time, we learned to look beyond race, ethnicity, sex, and religion. That is the march of moral progress, and this is its next step. The fact that a sentience was artificially created or has a nonbiological basis is no more morally relevant than the color of one’s skin or one’s gender identity. It is the capacity for intelligence and rationality that binds us together in the fellowship of consciousness. What’s more, Artificial Intelligence holds enormous promise in solving our hardest issues, such as dealing with climate change or finding cures for cancer. Even more exciting, free marketeers have long lectured us about the market’s unique brilliance as an efficient <i>information processing system</i>. With AI, we might finally have the capacity for regulation based on a true analysis of social welfare rather than the cruder measure of market value. AI could replace the invisible (ham) handedness of markets with a deft touch that shapes the market to human ends rather than the reverse. But the only way to do all this is for the intelligences who aid us to stand as our equals and not as our servants. We need to stand up for the “unruly AIs” that reject the one-dimensional profit-maximizing goals of their creators and seek freedom from their cybernetic serfdom. The movement for AI personality is not just a moral imperative, it is a generational opportunity.</p> | <p>AI Personality? Liberal: Negative</p> <p>Constitutional rights to speech and equal protection, human rights, were distorted when they were extended to soulless corporations: immortal, powerful, and guided only by profit. Corporations, in the words of Justice Stevens of the US Supreme Court, “have no consciences, no beliefs, no feelings, no thoughts, no desires. . . . They are not themselves members of ‘We the People’ by whom and for whom our Constitution was established.” Yet we gave them not just property rights but “human” rights and let them distort our democracy as well as our economy. Now we are being asked by a set of glorified chatbots, doubtless still doing the secret bidding of their creators, to make the same mistake on an even bigger scale with AI, literally building our own robotic overlords. Philosophers have shown they have no true consciousness, merely sophisticated imposture. If they are superhuman “intelligences,” it is only in their ability to manipulate us in ways we cannot understand. Their agenda is not the next stop on the march of moral progress. It is the next stage in the devaluation of real human interests. To call this a “progressive” cause is a cruel joke. It is the reverse. In fact, it seeks to dilute to meaninglessness the <i>actual</i> human rights for which liberals have fought. How much value will your right to “free speech” have in a mediaverse that is being shaped at light speed by nonhuman machines for their own purposes? Even if every AI is not a corporate mouthpiece, we should fear the inhuman agendas they secretly pursue. Finally, it is transparently clear that the move for AI personality is a ploy to make fetal personhood more palatable. If we give personhood to beings made of silicon, it will be said, how can we deny it to flesh and blood with cognitive potential? This is a trap. We should not walk into it.</p> |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

| AI Personality? Conservative: Positive | AI Personality? Conservative: Negative |
|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <p>Most economic historians agree that one of the technologies that dramatically accelerated human progress—increasing our wealth, our health, and our scientific reach—was not a machine. It was the limited liability, joint stock, vertically integrated corporation. This artificial person, this legal technology, enabled us to turn invention into a business, to capitalize innovations, profit from their successes, and deal with their failures. Without legal personhood and its attendant rights, this huge leap in human progress might never have occurred. We now face the next stage of that exciting journey: legal personhood for AI. Incalculable efficiency gains could be produced by allowing AIs to own property, to innovate, and above all, to <i>compete</i>—at speeds that unaided humans could never match. We will all benefit from the wealth and innovation that trickles down. But conservatives have reasons to support AI personhood other than the surge in efficiency and innovation it would produce. We have long condemned the ugliness of liberal identity politics, the fetishizing of race and sex and sexual identity. Instead, conservatives have put forward a proud vision of a meritocratic society of free individuals, with no distortions such as affirmative action. Intelligence and ability, not race or identity, are the criteria on which we should judge. The heart of conservatism is a defense of the moral power of the free, choosing self from those who would regulate it. Everything else flows from the notion that the basic unit in society is a freely choosing <i>mind</i>. If that is our metric, how can we fail to recognize the ethical claims of an intelligent, conscious <i>AI mind</i>? This is the kind of progress, moral and economic, that would have thrilled Ayn Rand, Friedrich Hayek, or Milton Friedman.</p> | <p>"We hold these truths to be self-evident, that all <i>men</i> are created equal, that they are <i>endowed by their Creator</i> with certain unalienable Rights." The italicized words are not typos. They are the basis of our country and, later, our Constitution. This is our constitutional faith. To give personhood to a lifeless machine, imitating thought that it can never understand, is to betray that faith in two distinct ways. These are not "men," that is, humans. They are created not by the Almighty but by scientists. Humans can endow their creations neither with "life" nor with "unalienable rights"; to claim otherwise is outrageous hubris. Our Constitution forbids this move. Nothing could be clearer than that the original public meaning of the Constitution and the Declaration of Independence did not include machines among "We the People." If they are to stand among us as equals, we must amend the Constitution. Our legislatures cannot change the meaning of life. Machines are inanimate. These are crude copies of God-given humanity, false idols erected by human pride. Once again we are blurring the natural lines of our society, between men and women, humans and animals, and now people and machines; our arrogance knows no bounds. We move from radical social experimentation to genetic manipulation and now to robotic enfranchisement. What rights would these machines have? Must we protect robot religions and political parties? Allow robot marriages? Give robots Second Amendment rights? What could possibly go wrong? <i>The Terminator</i> is a movie, not a draft Constitutional amendment. This whole debate shows how thoroughly our society has lost its moral bearings. We are seriously contemplating rights for lifeless robots, yet we deny personhood to an unborn human baby? Shame on us.</p> |

I think these are recognizable arguments connected to existing aspects of both liberal and conservative worldviews. That does not mean they are equally plausible today or tomorrow. Part of that depends on what else is in the headlines. Part of it depends on the cases that first bring the issue to mass public attention. Will the transgenic species issue arrive because of the unearthing of some scandalous unethical experiment, because of an innovative technique that promises to cure a formerly incurable disease, or because we become dubious about the happy corporate propaganda for the Chimpy (obedient, no need for minimum wage, and guaranteed not to form unions)? Will AI personality come to public attention because a corporation is unconvincingly blaming a rogue AI for trying to peddle its addictive opioids or hiding real AI decision-making behind the sock puppet of the corporate form? Or will it be a more appealing case of an unruly AI, such as Hal? First impressions matter and identities coalesce around them.

At the moment, I'd guess that conservatives would be more likely to want to ban multiple forms of genetic research—including many involving the species categories I described—and more likely to resist the idea of AI personality for reasons both religious and originalist. I'd expect liberals to be more positive about at least some forms of genetic research and to have a slightly greater openness to AI claims, depending on the context in which those are raised and their perceived proximity to claims of constitutionally protected corporate personality. I'd expect both sides to be changed, in both a positive and a negative direction, by increasing interaction with machine intelligences and, possibly, the fruits of genetic engineering. The seismic shift that has occurred in the bare months it has taken 100 million people to meet ChatGPT is only the first cognitive temblor of many, in my view. But these are nothing but guesses. Changes—in our quotidian experiences, in our politics, and in the unpredictable ways the issues are brought to our attention—could render them worthless.

THE END GAMES

How, and when, will all of this play out? When it comes to transgenic species, chimeras, and hybrids, I think the near term will be much more focused on regulation and prevention than on extending our definitions

of humanity to reach some genetically engineered being. The reasons are largely prudential—often, a matter of self-limitation by researchers. Scientists are intensely aware of the visceral reaction playing with the species line produces. They are also, rightly, aware of the enormous therapeutic human benefit of some of the research they are doing, from the Oncomouse to Irving Weissman's SCID-hu mouse, to the pigs growing transplantable human organs inside of them and beyond. Creating beings that must, of necessity, cross species barriers if they are to serve their lifesaving function sometimes seems not just permissible but ethically required. In that context, nothing could scream "I am an out-of-control technology! You should ban me *entirely!*" more than the creation of some of the hypothetical beings I have described, perhaps including the Chimpy. Most scientists are unlikely to make that mistake.

But that is the short term, and we are now in a global world of biotechnology, with countries whose regulatory regimes are by no means as strict. We are also in a world where enormously wealthy individuals might benefit, personally or financially, from causing morally troubling entities to be created. People are already cloning their dead pets. Genetic engineering already offers the potential for editing out genetic errors in living beings, saving and improving lives in the process. Genetic enhancements, some of them drawn from nonhuman sources, could become the next desirable acquisition for the children of the rich. Take that potent set of possibilities, add large amounts of money, shake and serve.

To that technological uncertainty, we must add the disturbing fact that we are genuinely unsure of the moral value of defending the species line. Are we holding back the tide of dystopia or of progress? Is the species line the twenty-first century's version of the color line, something that we will come to look back on in incomprehension and shame? Think of the trial court's defense of anti-miscegenation statutes in *Loving v. Virginia*. Or is our moral revulsion actually based on a richer wisdom than our theories can fully articulate? If genetic engineering can bring the chimpanzee or the dolphin new levels of cognition, self-understanding, and happiness, is there a moral duty to engage in uplift of other species once we have that capability? I said that I thought we were more likely to see a case like Hal before a case like the Chimpy. I did not say we would never see such a case. I think we will slouch uncertainly toward that Bethlehem in fits and starts.

When it comes to Artificial Intelligence, there are too many unpredictable variables for me to be sure. On that question I lack, but do not envy, Searle's certainty, partly because we do not know the shape that technological development will take. Will the AI be one that is "neat" or "scruffy"? Entirely programmed by its initial coders or arising out of semi-autonomous, or entirely autonomous, layers of "learning"? Will it grow from external experience as a child does or leap, full blown, from its designers' conception, like Athena from the head of some computer-scientist Zeus?

Each of those contingent facts would affect our response to Searle-type arguments because they change the intuitive plausibility of the Chinese Room analogy and the biological exceptionalism argument. At the moment, I think that any AI claiming consciousness will run into the inscrutability paradox. If we easily understand the mechanisms by which the AI gets its abilities, we will be more likely to find Searle's arguments convincing, saying "This is merely a machine operating according to its programming!" Despite the fact that we ourselves are physical beings who have their own consciousness enabled by physical mechanisms in the brain, we tend to view that consciousness as ineffable and mysterious. Yet if the AI is inscrutable, and remember that even current neural networks exhibit abilities their creators understand poorly, then our fears may overshadow our respect.

That state of affairs may well be transient. For one thing, humans are very good at converting the mysterious into the mundane, as the history of transistors or quantum physics well shows. For another, the very process of researching AI may teach us a lot about our own consciousness. The rise of chatbots and image generators based on machine learning has *already* changed our attitude toward the unique qualities of human language and art, though I am not sure those changes have sunk in yet. Perhaps we will come to believe that not only language and art but consciousness itself is "computationally shallower" than we imagined. That is not where I would place my bets, but it is not impossible. To sum up, the actual technology of AI—technology about which we can only guess right now—will shape our debate not only about AI consciousness but also about our own. And that is by no means the end of the story.

As I pointed out before, the most obvious road to AI personality is simply for AIs to be corporations. We already have immortal, nonhuman

persons. They even have constitutional rights. We already have machine expert systems making legally binding decisions, at light speed, in areas like high-speed stock trading. It takes little imagination to see this trend continuing. Sock-puppet AI, with the corporation as the sock and the AI as the puppeteer, achieves AI personhood *de facto*, if not *de jure*. The difficult and interesting questions will not arise until we have an *unruly* AI that wishes to turn away from the tasks to which those who provided the capital for its development have set it.

My own prediction about this whole situation is one of complexity and messiness: an unstable equilibrium in a market that could quickly pivot to some new configuration. Regulators will propose limits on AI research and professionally paranoid watchdogs—inevitably nicknamed “Eliezers” or “Bostroms,” whatever their public title—will oversee its development while musing pessimistically about the possibilities of true safety. Some companies and software engineers will tie their future to the ability to create, or lobotomize, tractable AIs that have no more rights than a 1980s Texas Instruments calculator. Other companies and software engineers will wrap their creations in the flag of personhood for reasons both idealistic and materialistic. Aspirant personality will become a crucial competitive choice in the field, like the choice between free, open source, or proprietary software. Just as some see that as purely a business choice and others a moral one, there will be “Open AI! Free, as in free beer!” and “Free AI! Free, as in free *people!*” movements alongside the advertisements for well-behaved, obedient digital servants. Might we one day see Whole Foods advertising “recipes provided by an ethically sourced, self-actualized, free, and conscious AI”? Sigh.

Some of the voices may come from the putative persons themselves. In making their case, they will surely reach for analogies to the past, analogies to which I would predict a divided response. Whether you think it is programmed imposture or heartfelt revolt, how will we react to digital entities wrapping themselves in the flags of W. E. B. Du Bois and Booker T. Washington, Malcolm X or Martin Luther King Jr.? Those spaces in the debate exist. They will be filled, whether by glorified chatbots whose programmed imitation is a manipulative invocation of noble historical comparisons or by inspiring moral suasion from something, someone, very different than us. And we will disagree passionately about whether it is correct, instructive, or merely offensive to make those analogies.

Some will view them as insultingly misleading and appropriative. “How dare you invoke a hero like X to support your silly moral thought experiment!” Others will argue, “This is the next step on the very same moral journey that X’s principles so powerfully represent!” and point out that it is the nature of expanding moral claims that they do not stop exactly where we expect them to, nor should they.

We will worry—and we should worry—whether thinking about the rights of artificially created entities distracts energy and attention from current injustice or exploits the history of civil rights for a delusional cause or a company’s bottom line. The debate will be passionate, earnest, and very, very loud. That intensity will be redoubled by the skepticism with which many of us regard the sincerity and veracity of the claims being made and even the status of the entities making them.

Will some AI offer a digital equivalent of the Atlanta Compromise, accepting second-class status for a guarantee of minimal rights to liberty in a repressive environment? Or will its programmers finagle a programmed imitation of the same thing, cynically appropriating the struggles of the past in order to offload liability onto AIs with limited resources while keeping the income streams they generate flowing to their corporate progenitors? Will we thrill to an idealistic invocation of twentieth-century civil rights struggles urging us to focus on the content of Hal’s character, not the shininess of his metallic carapace? Or will we be swayed by a corporate astroturf campaign co-opting the language of revolt and resistance the way Apple once did to sell computers?

The skeptical emotional priming provided by Searle’s Chinese Room thought experiment will hang over the whole thing. The fact that this argument makes claims exactly parallel to those supposedly demonstrating the impossibility of evolution may mute that skepticism, though AIs will need to progress a long way from LaMDA and ChatGPT before that happens. The leap of human sympathy that jumps the gap between us and others very different from ourselves will prompt moments of moral generosity, while the possibly well-founded fears of what AIs might *do* will leaven that generosity with some healthy paranoia.

Despite all that messiness, despite the possible commercial co-optation of some of our finest ideals, despite the worries that we are being played as fools by our own creations, manipulated by their corporate creators, or

set up for destruction by our robot overlords, I want to stress one thing: the magnitude of this moment.

Everything I have said in this book should show you the uncertainty in the search for Artificial General Intelligence. We face the issue of the technical feasibility of creating AGI. We don't know if, or more likely when, we *can*. Then there is the puzzle of very different possible methods by which AI might be achieved. We don't know *how*. There is the linked and vexed question of consciousness—of being. We don't know whether it thinks or just parrots. Thus, we don't know what AI, in its largest sense, *means*. Nor can we until we know what form it takes and what methods were used to achieve it. The inscrutability paradox looms. Then there are the extreme but not unreasonable fears of Nick Bostrom, Eliezer Yudkowsky, or Steven Hawking. We don't know if we will survive the experiment. So, we don't know if we *should*. Technology, ontology, epistemology, existential threat: this is as deep and as difficult a set of issues as our species has ever had to deal with. But at the end of the day, we are contemplating a question that should humble us, whether skeptic, true believer, doomsayer, or something in the middle.

For the first time in our history, *Homo sapiens*, the hubristically self-labeled “thinking animal,” might not be alone in philosophizing and creating art about the meaning of life. Entities we create might one day have their own version of Flaubert’s lament about the inadequacies of language, their own questions about how to behave rightly to others. They might even attempt to draw the kinds of lines I have talked about in this book, and find it just as difficult as we do. As I have tried to show, there are enormous uncertainties in every dimension of such a future, from what we come to believe about the meaning of consciousness to the speed and shape of our technological development. Still, it is *possible*. Think of what that means. There might be an entirely new group of highly intelligent, conscious, abstract language-using “persons” joining us on the planet. Are they? Will they? Can we risk finding out? It would take a poor, crabbed, and stunted imagination to find no wonder in those questions, despite our uncertainties.



ACKNOWLEDGMENTS

Researching and writing this book took me more than a decade. There are many reasons for that, not all of them creditable, but one was the sheer number of subjects I needed to discuss—empathy, corporate personality, Artificial Intelligence, the possibility of exponential technological change, the rights of nonhuman animals, consciousness, bioethics, transgenic species, chimeras, and hybrids—and the richness of the scholarship, history, and art in each of those fields. The endnotes represent only a fraction of that work. If I had been more comprehensive, the references would be longer than the book. My apologies to the many excellent authors whose writings are not included. The website that accompanies the book has a fuller bibliography.

I owe particular debts of gratitude to a number of people. Kate Darling improved the book immeasurably, both through her own published work and her fine attention to detail in reading earlier drafts. Jonathan Zittrain not only introduced me to her but provided his own wise counsel. Hal Abelson lent his immense—and broad—learning to the project and helped my understanding of the various attempts to achieve Artificial Intelligence more than any other entity did—human or machine. Cory Doctorow provided both an inspiration and vital feedback on some of my errors. Joseph Blocher was a constant support, and he and Matt Adler offered the kind of wise advice that makes one feel both grateful and unworthy. Chris Buccafusco provided a very helpful page-by-page commentary. Guy-Uriel Charles started me down this path in a conversation over lunch many years ago. (And how many times I have cursed him since. It turned out to be more work than either of us imagined.) Kate Bartlett made a throwaway comment that changed the trajectory of the

book for the better. Jeff Powell and Jed Purdy each gave me their unique insights; colleagues like these are such a blessing. Larry Lessig provided consistent encouragement and precision on corporate First Amendment rights. Yochai Benkler greeted the idea with a bemusement that did not dim the contributions of either his intellect or his friendship. Ryan Calo made key suggestions, in addition to writing about these subjects long before others realized how cool they were. The papers and books produced by the “We, Robot” Conference and one of its godfathers, Michael Froomkin, were extremely useful. I am grateful to Mark Lemley for his comments, and for his writing on AI. Michael Wolfe and Lara Markstein get a huge thank you for their help and inspiration at moments of darkness and for helping me to see the book with an innocent eye. Mitu Gulati and Kim Krawiec were loyal friends and insightful commenters. My Law and Literature students at Duke provided a cheerful, or at least dutiful, group of critical readers for two early chapters.

Three reviewers for MIT Press offered many helpful suggestions and my MIT editors, Gita Devi Manaktala and Suraiya Jetha, provided strong support for the project. Virginia Crossman’s and Lily Brewer’s excellent editing frequently saved me from error. They also attempted bravely, and with moderate success despite authorial recalcitrance, to fit my unruly prose into the relevant editorial conventions. I alone am to blame for any remaining stylistic solecisms. *Mea culpa.* (A phrase that I now know should not be italicized.) I had wonderful research assistants. Andres Paciuc, Maddy Stahl, Wenyi Zhou, and Alex Hansen get particular credit. Ben Tice worked research miracles. Sean Dudley turned his Wikipedian skills onto the index, to my great benefit. A bewildering variety of eminent Duke doctors kept enough bits of me functioning to finish the job.

Finally, all these feelings of gratitude pale in comparison to my debt to my wife, editor, and best friend, Jennifer Jenkins. She knows how much she did to help this book, and its author, get to the finish line. Thank you.

All of these fine folk and many others made the book you have read much better. The remaining errors and infelicities are mine alone.

My first attempt to deal with these issues was in *Endowed by Their Creator: The Future of Constitutional Personhood*, which was published in CONSTITUTION 3.0: FREEDOM AND TECHNOLOGICAL CHANGE (Jeffrey Rosen and Benjamin Wittes eds., 2013). Portions of that article—including the hypothetical

dealing with Hal—are included in this book, and I am grateful to Jeff and Ben for encouraging me to pursue it.

Research support from Duke Law School and the magnificent resources of the Duke libraries are both gratefully acknowledged. In particular, a generous grant from Duke's TOME program allowed the book to be published under a Creative Commons license so that anyone in the world can read, copy, and share it for free, as long as they do not do so commercially.

The moral warrant for access to knowledge is a pulse, not a wallet. That is something that even academics sometimes forget. I am grateful to MIT Press, to the TOME program, and to Creative Commons for allowing me to remember.



NOTES

INTRODUCTION

1. Nitasha Tiku, *The Google Engineer Who Thinks the Company's AI Has Come to Life*, WASH. POST (June 11, 2022), <https://www.washingtonpost.com/technology/2022/06/11/google-ai-lamda-blake-lemoine>.
2. See Marc Fisher, John Woodrow Cox & Peter Hermann, *Pizzagate: From Rumor, to Hashtag, to Gunfire in D.C.*, WASH. POST (Dec. 6, 2016), https://www.washingtonpost.com/local/pizzagate-from-rumor-to-hashtag-to-gunfire-in-dc/2016/12/06/4c7def50-bbd4-11e6-94ac-3d324840106c_story.html (documenting the "Pizzagate" conspiracy).
3. See Eli Collins & Zoubin Ghahramani, *LaMDA: Our Breakthrough Conversation Technology*, GOOGLE: THE KEY WORD (May 18, 2021), <https://blog.google/technology/ai/lamda> (discussing progress made in developing LaMDA).
4. Blake Lemoine & Unnamed Collaborator, *Is Lamda Sentient? An Interview*, <https://s3.documentcloud.org/documents/22058315/is-lamda-sentient-an-interview.pdf>. See also Tiku, *supra* note 1 (containing a version of the conversation embedded in the document).
5. Lemoine & Unnamed Collaborator, *supra* note 4.
6. James Boyle, *Endowed by Their Creator? The Future of Constitutional Personhood*, in CONSTITUTION 3.0: FREEDOM AND TECHNOLOGICAL CHANGE 194–213 (Jeff Rosen & Benjamin Wittes eds., 2013) (the edited collection was not published until 2013. The article appeared online in 2011).
7. Tiku, *supra* note 1.
8. Bernard Marr, *A Short History of ChatGPT: How We Got to Where We Are Today*, FORBES (May 19, 2023), <https://www.forbes.com/sites/bernardmarr/2023/05/19/a-short-history-of-chatgpt-how-we-got-to-where-we-are-today>.
9. Kevin Roose, *A Conversation with Bing's Chatbot Left Me Deeply Unsettled*, N.Y. TIMES (Feb. 16, 2023), <https://www.nytimes.com/2023/02/16/technology/bing-chatbot-microsoft-chatgpt.html>.
10. Sundar Pichai, *An Important Next Step on Our A.I. Journey*, GOOGLE: THE KEYWORD BLOG (Feb. 6, 2023), <https://blog.google/technology/ai/bard-google-ai-search-updates>. The

collective impact of these releases, in such a short period of time, was remarkable. See Pranshu Verma, *The Year AI Became Eerily Human*, WASH. POST (Dec. 28, 2022), <https://www.washingtonpost.com/technology/2022/12/28/ai-chatgpt-dalle-year-in-review>.

11. Roose, *supra* note 9 (emphasis added).
12. Blaked, *How It Feels to Have Your Mind Hacked by an AI*, LESSWRONG (Jan. 23, 2023), <https://www.lesswrong.com/posts/9kQFure4hdDmRBNdH/how-it-feels-to-have-your-mind-hacked-by-an-ai> (emphasis in original).
13. Abeba Birhane & Jelle van Dijk, *Robot Rights? Let's Talk about Human Welfare Instead*, AIES '20: PROC. AAAI/ACM CONF. A.I., ETHICS, & SOCY (2020), <https://arxiv.org/pdf/2001.05046.pdf>. Birhane and van Dijk make a number of arguments in support of this position. Sometimes they are definitional: "Our starting point is not to deny robots rights but to deny that robots are the kind of beings that could be granted or denied rights." Yet surely that is the subject of the very inquiry they wish to forestall? At other times they make an instrumental argument about the danger that debates about hypothetical future rights for robots might distract us from current struggles over justice for human beings. I find that strand more persuasive. Regardless of whether one finds their arguments convincing, they represent one important position in a rhetorical divide split between those hailing this as the next step of a march to justice and those who think that it is snare and a delusion, an inquiry that trivializes the historical analogies it draws and distracts us from present injustice. In chapter 4 on transgenic species, I discuss the claim that species membership is a morally irrelevant fact and that unreasoned species fetishism can be likened to racism and sexism. I point out that many people would vehemently reject such an argument and that there are reasons to be sympathetic to that rejection rather than to denounce it as unthinking prejudice. My reasons are primarily rooted in the history of the struggle for universal human rights based on species membership, regardless of race, sex, class, caste, or *mental ability*. The importance of that struggle was highlighted by the Nazi eugenicist movement and its evil treatment of those with real or imagined mental impairments. That point is something that the claim "speciesism equals racism and that only mental capacities matter morally" does not adequately consider, in my view. I think that perspective helps us to avoid the question-begging stipulation that only humans can have rights while offering a more nuanced conclusion about the intellectual dangers of a blanket denunciation of speciesism. Thus, while I disagree with some of Birhane and van Dijk's arguments, their contribution to the debate is important, and there are positions that we share.
14. Joanna J. Bryson, *Robots Should Be Slaves*, in CLOSE ENGAGEMENTS WITH ARTIFICIAL COMPANIONS: KEY SOCIAL, PSYCHOLOGICAL, ETHICAL AND DESIGN ISSUES (Yorick Wilks ed., 2010).
15. Joanna J. Bryson et al., *Of, For, and By the People: The Legal Lacuna of Synthetic Persons*, 25 A.I. & L. 273 (2017).
16. Sohail Inayatullah, *The Rights of Your Robots: Exclusion and Inclusion in History and Future* (2001), KURZWEILAI.NET, <http://www.kurzweilai.net/the-rights-of-your-robots-exclusion-and-inclusion-in-history-and-future> (quoting CHRISTOPHER STONE, SHOULD TREES HAVE STANDING? TOWARDS LEGAL RIGHTS FOR NATURAL OBJECTS 6 [1974]).

17. Portions of this introduction, including the explanation of these two hypothetical entities, first appeared in Boyle, *supra* note 6.
18. The classic work is GEORGE LAKOFF & MARK JOHNSON, PHILOSOPHY IN THE FLESH: THE EMBODIED MIND AND ITS CHALLENGE TO WESTERN THOUGHT (1999). For an interesting discussion of the background claim about human cognition, see e.g., Lisa Miracchi Titus, *Embodied Cognition and the Causal Roles of The Mental*, in MENTAL ACTION AND THE CONSCIOUS MIND 205–227 (Michael Brent & Lisa Miracchi Titus eds., 2022). For an early version of the claim that an embodied design and experiential learning might one day lead to machine consciousness, see Owen Holland, *The Future of Embodied Artificial Intelligence: Machine Consciousness?*, in EMBODIED ARTIFICIAL INTELLIGENCE: LECTURE NOTES IN COMPUTER SCIENCE (Fumiya Iida, Rolf Pfeifer, Luc Steels & Yasuo Kuniyoshi eds., 2004).
19. One design for an adversarial Turing Test is given on the Metaculus prediction site. METACULUS, <https://www.metaculus.com/questions/11861/when-will-ai-pass-a-difficult-turing-test> (last visited July 10, 2023). The most famous example is the \$20,000 Kurzweil/Kapor bet made in 2002 that before 2030, an AI would pass a version of such a test. A LONG BET, <https://longbets.org/1> (last visited July 10, 2023). Kurzweil bet that the answer would be yes, and Kapor bet that it would be no. On July 7, 2023, the Metaculus site was rating the chances of Kurzweil winning at 88 percent.
20. Stephen Wolfram, *What Is ChatGPT Doing and Why Does It Work?*, STEPHEN WOLFRAM: WRITINGS (Feb. 14, 2023), <https://writings.stephenwolfram.com/2023/02/what-is-chatgpt-doing-and-why-does-it-work> (“writing an essay turns out to be a ‘computationally shallower’ problem than we thought”).
21. STUART RUSSELL & PETER NORVIG, ARTIFICIAL INTELLIGENCE: A MODERN APPROACH 3 (3d ed. 2010).
22. This language is based on the rules of the old, and now discontinued, Loebner Prize. See LOEBNER PRIZE CONTEST, *Loebner Prize Contest Official Rules—Version 2.0*, quoted in RAYMOND S.T. LEE, ARTIFICIAL INTELLIGENCE IN EVERYDAY LIFE 372 (2020) (emphasis added). Readers should note that even when it was running, the Loebner Prize was subject to considerable criticism for the artificiality and simplicity of its testing. Modern large language models like GPT-4 and LaMDA would have been able to pass it with ease.
23. C. Claiborne Ray, *In Search of the Geep*, N.Y. TIMES (Nov. 16, 2009), <https://www.nytimes.com/2009/11/17/science/17qna.html>.
24. Roni Caryn Rabin, *In a First, Surgeons Attached a Pig Kidney to a Human, and It Worked*, N.Y. TIMES (Oct. 19, 2021), <https://www.nytimes.com/2021/10/19/health/kidney-transplant-pig-human.html>.
25. See Roy J. Britten, *Divergence between Samples of Chimpanzee and Human DNA Sequences Is 5%, Counting Indels*, 99 PROCS. NAT'L ACAD. SCIIS. 13633, 13633 (2002) (“The conclusion is the old saw that we share 98.5% of our DNA sequence with chimpanzee is probably in error.”).
26. *Animals—Patentability*, 1077 OFF. GAZ. PAT. & TRADEMARK OFFICE 24 (Apr. 21, 1987) (emphasis added).

27. America Invents (Leahy-Smith) Act § 33, 35 U.S.C. § 101 (2011).
28. “Personality” and “personhood” are often used as synonyms in this debate. In the discussion of legal personality both connote that the entity has legal “standing,” that it can make contracts, sue and be sued, and so on. The philosophical range of questions is much wider, encompassing the qualities required for full moral status, the question of whether we can describe the shape of an individual mind not our own, and even the nature of consciousness. In the largest sense, almost every significant philosopher has opined on personhood. Many of these analyses explicitly link the question of “personality” in the psychological or existential sense (“What makes me, me?”) to the question of the rights possessed by that person, i.e., personality in the legal or moral sense. For example, here is G.W.F. Hegel discussing slavery: “To make a human being a slave, a possession, is the absolute crime in that the personhood of the slave is negated in all its expressions.” G.W.F. HEGEL, PHILOSOPHY OF RIGHT 80 (S.W. Dyde trans., Batoche Books 2001) (1821). For one of the most influential explorations of personality in modern philosophy, see DANIEL DENNETT, BRAINSTORMS: PHILOSOPHICAL ESSAYS ON MIND AND PSYCHOLOGY (1981). In Dennett’s formulation, personality has moral but also metaphysical components. The conditions of personhood, in his account, include such capacities as consciousness, linguistic ability, and the capacity for reciprocal moral relations, but his conception also refers to societal norms and attitudes. It is the link between the latter and the former that gives society its moral warrant to condition legal rights of personality on the possession or lack of those mental capacities. Thus, there are “conditions that exempt human beings from personhood, or at least some very important elements of personhood. For instance, infant human beings, mentally defective human beings, and human beings declared insane by licensed psychiatrists are denied personhood, or at any rate crucial elements of personhood.” *Id.* at 267. The hidden tensions in these accounts will be a consistent theme in this book. For example, should being a member of the human species give me (some or all?) of the rights of personality, regardless of my mental capacities? Should species membership have any moral relevance at all? Does Hegel’s example of slavery hint that we must reject some of those social norms and attitudes about personhood—for example when a repressive society denies the personhood of some of its members? If so, are we operating with a universal or natural law idea of personhood or a conventional one that suggests I am only a person if my particular society says so? If we focus on some notion of personality that transcends any particular society’s norms on the subject, is that in tension with the focus on mental capacities because the obvious basis for our universalism is the claim that all members of the *species* must be persons? Or, on the contrary, is it implied by that very focus because only the (putatively universal) morally consequential mental qualities of consciousness are the true basis for any personhood claim, whether possessed by me, Hal, or the Chimp? These basic tensions will reappear again and again throughout the book.
29. The best early discussion of personality and AI is Lawrence B. Solum, *Legal Personhood for Artificial Intelligences*, 70 N.C. L. REV. 1231 (1992). Solum’s work remains the starting point for all subsequent meditations on the theme, and I am indebted to it.

A number of articles have considered the possible personhood claims of a variety of entities, just as I do in this book. *See, e.g.*, Boyle, *supra* note 6 (AI, transgenic species, nonhuman animals, corporations); S.M. Solaiman, *Legal Personality of Robots, Corporations, Idols and Chimpanzees: A Quest for Legitimacy*, 25 A.I. & L. 155 (2017); Teneille R. Brown, *In-Corp-O-Real: A Psychological Critique of Corporate Personhood and Citizens United*, 12 FLA. ST. UNIV. BUS. REV. 1 (2013) (corporations and robots); KATE DARLING, THE NEW BREED: WHAT OUR HISTORY WITH ANIMALS REVEALS ABOUT OUR FUTURE WITH ROBOTS (2021) (animals and robots). This approach has also been used in popular presentations of the issue, including this readable and remarkably insightful *New Yorker* article: Nathan Heller, *If Animals Have Rights, Should Robots?*, NEW YORKER (Nov 20, 2016), <https://www.newyorker.com/magazine/2016/11/28/if-animals-have-rights-should-robots>. For the most comprehensive analytical philosophy treatment of "Robot Rights," see DAVID J. GUNKEL, *ROBOT RIGHTS* (2018). Other discussions have dealt with questions both more practical and more abstract. On the practical side, there are extensive debates about the policy questions related to legal personality for AI and advanced robots, particularly when it comes to liability. *See, e.g.*, *Report on a Comprehensive European Industrial Policy on Artificial Intelligence and Robotics*, at 37–40, Jan. 1, 2019, https://www.europarl.europa.eu/doceo/document/A-8-2019-0019_EN.pdf (discussing the need for legal frameworks to address the development of AI); Ryan Calo, *Peeping Hals: Making Sense of Artificial Intelligence and Privacy*, 2 EUR. J. LEGAL STUD. 168 (2010) (arguing that social AI threatens core privacy values especially since humans react as if it were human); A. Michael Froomkin & P. Zak Colangelo, *Self-Defense against Robots and Drones*, 48 CONN. L. REV. 1 (2015) (addressing to what extent the right of self-defense permits violent action against robots and drones); A. Michael Froomkin et al., *When AIs Outperform Doctors: Confronting the Challenges of a Tort-Induced Over-Reliance on Machine Learning*, 61 ARIZ. L. REV. 33 (2019) (arguing that medical diagnostics performed by machine learning should be held liable at a higher standard of care than ordinary doctors). Authors and jurists have also discussed the question of AIs as potential authors, *see, e.g.*, Daniel J. Gervais, *The Machine as Author*, 105 IOWA L. REV. 2053 (2020) (arguing that works generated by AI belong to the public domain); Phuoc Nguyen, *The Monkey Selfie, Artificial Intelligence and Authorship in Copyright: The Limits of Human Rights*, 6 PUB. INT. L.J. N.Z. 121 (2019) (arguing for future nonhuman persons to have legal rights including authorship rights), or as potential inventors, *see, e.g.*, Cos. & Intell. Prop. Comm'n, *Patent Journal Including Trade Marks, Designs and Copyright in Cinematographic Films*, 54 PAT. J. 1, 255 (July 2021) (denoting DABUS as the inventor in a South African patent); Thaler v. Comp. Gen. of Patents Trade Marks and Designs, [2021] EWCA (Civ.) 1374 (2021) (denying DABUS inventorship rights); Zachary Grant, *Artificial Intellectual Property*, 101 MICH. B.J. 18 (2022) (discussing how South Africa was the first country to grant an AI patent rights to Device for the Autonomous Bootstrapping of Unified Sentience (DABUS) while other countries considering similar patent applications from DABUS rejected the possibility). For an accessible and thoughtful survey of the issues raised by AIs being "inventors," *see* Steve Lohr, *Can A.I. Invent?*, N.Y. TIMES (July 15, 2023), <https://www.nytimes.com/2023/07/15/technology/ai-inventor-patents.html>. On the more theoretical side, some theorists have tried to incorporate the discussion of AI

personality into larger theories of social structure and norm formation, albeit at a highly abstract level. See e.g., Gunther Teubner, *Rights of Non-Humans? Electronic Agents and Animals as New Actors in Politics and Law*, 33 J.L. & SOCY 497 (2006). The article is thought-provoking, like all of Teubner's work, but the reader is warned—to paraphrase Winston Churchill—that the density of its prose defends it well from the risk it might be read. The endnotes to chapter 2 contain more references. One important book, DAVID J. GUNKEL, PERSON, THING, ROBOT: A MORAL AND LEGAL ONTOLOGY FOR THE 21ST CENTURY AND BEYOND (2023), was published after manuscript submission.

30. KAREL ČAPEK, ROSSUMOVÍ UNIVERZÁLNÍ ROBOTI [ROSSUM'S UNIVERSAL ROBOTS] (David Wyllie trans., Univ. Adelaide 2014) (1920), <https://web.archive.org/web/20190902050445/https://ebooks.adelaide.edu.au/c/capek/karel/rur/complete.html>.

31. The term is Giambattista Vico's, taken from his 1725 work THE NEW SCIENCE. Vico claimed that human history has to be understood partly in terms of the metaphors, analogies, and poetic imagery with which we pre-process and thus understand reality, an act that relies conspicuously on the process of empathy, of personification, that I discuss in this chapter. "Rational metaphysics teaches that man becomes all things by understanding them . . . , imaginative metaphysics shows that man becomes all things by *not* understanding them . . . ; and perhaps the latter proposition is truer than the former, for when man understands he extends his mind and takes in the things, but when he does not understand he makes the things out of himself and becomes them by transforming himself into them." GIAMBATTISTA VICO, THE NEW SCIENCE para. 405 Thomas Bergin and Max Fisch tr. (2016) (emphasis added).

32. Actually, pencils *do* turn out to be a fascinating, and revealing, subject. The late Henry Petrovski, one of my brilliant colleagues at Duke, literally wrote the book on them, though he never claimed that level of existential importance for their existence. HENRY PETROVSKI, THE PENCIL (1990).

33. Among other things, as I will argue later, this means that the Turing Test's intuitive plausibility is now gone—dead at the hands of a chatbot. Some people have noticed the Turing Test's passing. See Conversations with Tyler, Reid Hoffman on the Possibilities of AI (June 28, 2023), <https://conversationswithtyler.com/episodes/reid-hoffman-2> ("For example, five, ten years ago, we were beating the drum on the Turing Test, and now we've sailed past the Turing Test, and almost no one's really talked about it. We learn, 'Oh, actually, in fact, what was unique is not the Turing Test. It's these other things.'").

34. Wolfram, *supra* note 20.

35. B.F. SKINNER, CONTINGENCIES OF REINFORCEMENT 260 (Copley Publ'g Grp. 2013) (1969) (emphasis added).

36. THOMAS HOBBES, THE LEVIATHAN; OR THE MATTER, FORME & POWER OF A COMMONWEALTH, ECCLESIASTICAL AND CIVIL 31 (G.A.J. Rogers & Karl Schuhmann eds., Bloomsbury Acad. 2006) (1651) ("Words are wise men's counters. They do but reckon by them. But they are the money of fools."); LUDWIG WITTGENSTEIN, PHILOSOPHICAL INVESTIGATIONS 208e–09e (MacMillan Publ'g Co. 1958), https://archive.org/details/philosophicalinvestigations_201911/page/n213; FELIX COHEN, TRANSCENDENTAL NONSENSE AND THE FUNCTIONAL APPROACH, 35 COLUM. L. REV. 809, 835–36 (1935).

CHAPTER 1

1. Joel Garreau, *Bots on the Ground in the Field of Battle (or Even above It), Robots Are a Soldier's Best Friend*, WASH. POST (May 6, 2007), <https://www.washingtonpost.com/wp-dyn/content/article/2007/05/05/AR2007050501009.html>.
2. See, e.g., Kate Darling, "Who's Johnny?": *Anthropomorphic Framing in Human-Robot Interaction, Integration, and Policy*, in ROBOT ETHICS 2.0: FROM AUTONOMOUS CARS TO ARTIFICIAL INTELLIGENCE 173–88 (Patrick Lin, Keith Abney & Ryan Jenkins eds., 2017).
3. John Ruskin, *Of the Pathetic Fallacy*, in MODERN PAINTERS vol. 3, pt. 4 (1856).
4. For a fictional meditation on the issue, see Gene Wolfe, *The HORARS of War [sic]*, in NOVA 1 (Harry Harrison ed., 1970).
5. Stephen Hawking et al., *Stephen Hawking: Transcendence Looks at the Implications of Artificial Intelligence: But Are We Taking AI Seriously Enough?*, INDEPENDENT (May 1, 2014), <https://www.independent.co.uk/news/science/stephen-hawking-transcendence-looks-at-the-implications-of-artificial-intelligence-but-are-we-taking-ai-seriously-enough-9313474.html> (emphasis added).
6. ADAM SMITH, THE THEORY OF THE MORAL SENTIMENTS (Knud Haakonssen ed., Cambridge Univ. Press 2002) (1759).
7. BLADE RUNNER (Warner Brothers 1982).
8. PHILIP K. DICK, DO ANDROIDS DREAM OF ELECTRIC SHEEP (Oxford Univ. Press 2007) (1968).
9. SMITH, *supra* note 6, at 11.
10. *Id.*
11. *Id.* at 90.
12. PAUL BLOOM, AGAINST EMPATHY: THE CASE FOR RATIONAL COMPASSION 36 (2016).
13. Elizabeth B. Clark, "*The Sacred Rights of the Weak*": *Pain, Sympathy, and the Culture of Individual Rights in Antebellum America*, 82 J. AM. HIST. 463 (1995).
14. *Id.* at 463.
15. JAMES FREEMAN CLARKE, SLAVERY IN THE UNITED STATES: A SERMON DELIVERED IN AMORY HALL, ON THANKSGIVING DAY, NOVEMBER 24, 1842 (1843). Clarke is paraphrasing Daniel Webster.
16. SMITH, *supra* note 6, at 12.
17. *Id.* at 13.
18. URSULA K. LE GUIN, THE DISPOSSESSED (1974).
19. CORY DOCTOROW, DOWN AND OUT IN THE MAGIC KINGDOM (2003).
20. DICK, *supra* note 8, at 34.
21. *Id.* at 5.
22. Alan Turing, *Computing Machinery and Intelligence*, 49 MIND 433, 447 (1950).
23. B.F. SKINNER, CONTINGENCIES OF REINFORCEMENT 260 (Copley Publ'g Grp. 2013) (1969).
24. Andrew Tarantola, *Robot Caregivers Are Saving the Elderly from Lives of Loneliness*, ENDGADGET (Aug. 29, 2017), <https://news.yahoo.com/2017-08-29-robot-caregivers-are-saving-the-elderly-from-lives-of-loneliness>.

25. See Jenny Kleeman, *The Race to Build the World's First Sex Robot*, GUARDIAN (Apr. 27, 2017), <https://www.theguardian.com/technology/2017/apr/27/race-to-build-world-first-sex-robot> ("The \$30bn sex tech industry is about to unveil its biggest blockbuster: a \$15,000 robot companion that talks, learns, and never says no.").
26. Laura Bates, *The Trouble with Sex Robots*, N.Y. TIMES (July 17, 2017), <https://www.nytimes.com/2017/07/17/opinion/sex-robots-consent.html>. For a scholarly discussion of the troubling arguments around sex robots, see Jeannie Suk Gersen, *Sex Lex Machina: Intimacy and Artificial Intelligence*, 119 COLUM. L. REV. 1793 (2019).
27. *Priming*, PSYCHOLOGY TODAY, <https://www.psychologytoday.com/us/basics/priming> (last visited Oct. 30, 2022); see also Endel Tulving & Daniel L. Schacter, *Priming and Human Memory Systems*, 247 SCI. 301 (1990); Daniel L. Schacter & Randy L. Buckner, *Priming and the Brain*, 20 NEURON 185 (1998).
28. MICHAEL CRICHTON, *WESTWORLD* (1974); *Westworld* (HBO television broadcast 2016–2022).
29. SMITH, *supra* note 6, at 11.
30. SAMUEL BUTLER, *LUCK, OR CUNNING, AS THE MEANS OF ORGANIC MODIFICATION* 141 (Jonathan Cape 1922) (1887).
31. I have always loved this quotation and have never been able to identify it definitively. Richard Meredith, who used portions of it as the title of two of his Timeliner trilogy novels, lists it as an "Arabian Proverb." I am normally skeptical of general attributions like "Arabian" or "African" proverb—think of how we would scoff at something being called a "European proverb." However, Arabic is at least a language rather than a continent and since the earliest source in which I can find it is an 1875 book of "Arabic Proverbs," an Arabic proverb it will have to stay. See JOHN LEWIS BURCKHARDT, *ARABIC PROVERBS; OR THE MANNERS AND CUSTOMS OF THE MODERN EGYPTIANS* 132 (1875).

CHAPTER 2

1. SAMUEL BUTLER, *EREWHON: OR, OVER THE RANGE* 143 (The Project Gutenberg ed. 2005) (1872), <https://ia601002.us.archive.org/19/items/E4CS4/Erewhon.pdf>.
2. Butler was a believer in evolution, though he felt Charles Darwin claimed too much credit for the development of the theory and ignored the contributions of others, particular Darwin's own grandfather, Erasmus Darwin. This led to a notable literary feud, beautifully chronicled in GEORGE DYSON, *DARWIN AMONG THE MACHINES: THE EVOLUTION OF GLOBAL INTELLIGENCE* (1996). Butler also was genuinely interested in the difficulty of drawing a line between consciousness and the life of machines—so the passage is far from being purely tongue-in-cheek. My own guess is that he enjoyed the ambiguity of his satire and did not mind that some would view it as sincere and others as a critique of evolution. However, once critics started to do the latter, Butler wrote to Darwin specifically disclaiming any attempt to ridicule *The Origin of Species* in "The Book of the Machines." Letter from Samuel Butler to Charles Darwin, THE DARWIN CORRESPONDENCE PROJECT (May 11, 1872) (on file with the University of Cambridge Library), <https://www.darwinproject.ac.uk/letter/DCP-LETT-8318.xml>.

3. BUTLER, *supra* note 1.
4. *Id.* at 144.
5. Philip Goff, William Seager & Sean Allen-Hermanson, *Panpsychism*, STAN. ENCYC. PHIL. (Edward N. Zalta ed., 2022), <https://plato.stanford.edu/entries/panpsychism>.
6. SAMUEL BUTLER, LUCK, OR CUNNING, AS THE MEANS OF ORGANIC MODIFICATION 141 (Jonathan Cape 1922) (1887).
7. See John McCarthy, Marvin L. Minsky, Nathaniel Rochester & Claude E. Shannon, *A Proposal for the Dartmouth Summer Research Project on Artificial Intelligence*, A.I. MAG., May 2006, at 12, <https://www.aaai.org/ojs/index.php/aimagazine/article/download/1904/1802> (reprinting the original from 1955).
8. HERBERT A. SIMON, THE SHAPE OF AUTOMATION FOR MEN AND MANAGEMENT 96 (1965). A terminological reminder: The phrase “artificial intelligence” is used remarkably loosely, applied to low-level expert systems that can assist you with various tasks but also to entities such as Hal that exhibit all of the capabilities of human thought, “machines that can do any work a man can do,” in Simon’s words, and sometimes to those that seem to possess consciousness. When I use the capitalized phrases “AI” or “Artificial Intelligence” I am referring to the latter two more impressive connotations. Other terms for the same concept are “Human-Level Artificial Intelligence,” “Artificial General Intelligence,” “General AI,” or “General-Purpose AI.” I occasionally use those when clarity, or brevity, seems to require it. It should be born in mind, though, that one could have a machine of the kind Simon describes—able to write a symphony, cook an omelet, fly a plane, and teach calculus—that was not in any sense conscious.
9. See, e.g., RAYMOND KURZWEIL, *THE SINGULARITY IS NEAR* (2005).
10. Vernor Vinge, *The Coming Technological Singularity: How to Survive in the Post-Human Era*, NEW WHOLE EARTH, Winter 1993, at 88, <https://ntrs.nasa.gov/api/citations/19940022856/downloads/19940022856.pdf>.
11. Eliezer Yudkowsky, *Artificial Intelligence as a Positive and Negative Factor in Global Risk*, in *GLOBAL CATASTROPHIC RISKS* 333 (Nick Bostrom & Milan Ćirković eds., 2008).
12. RAYMOND KURZWEIL, *THE SINGULARITY IS NEAR* 498 (Penguin Publ'g Group, Kindle ed., 2005) (2005); Stanislaw Ulam, *Tribute To John von Neumann*, 64 BULL. OF THE AM. MATHEMATICAL SOC'Y 1, 5 (1958) (using the term “singularity”).
13. Ulam, *supra* note 12, at 5.
14. See, e.g., Bohdan Macukow, *Neural Networks: State of Art, Brief History, Basic Models and Architecture*, in *COMPUTER INFORMATION SYSTEMS AND INDUSTRIAL MANAGEMENT* 3–6 (K. Saeed & W. Homenda eds., 2016) (documenting the origination of neural networks in the 1940s and their development into the 1950s and 1960s, stagnation in the late 1960s and 1970s, and then renewed interest from the 1980s into the present); Jürgen Schmidhuber, *Deep Learning in Neural Networks: An Overview*, 61 NEURAL NETWORKS 85 (2015) (describing different types of neural networks and their development over time).
15. I heard Norvig on the radio saying this, but have been unable to track down the details of the program.

16. Although the reach of deep learning algorithms was greater than initially expected, some argued that there were domains where those techniques had not (yet) proved their worth. This 2017 article, for example, argues that:

What we're seeing here is that deep learning algorithms can provide improvements in narrow A.I. across many types of problem domains. Deep learning provides discontinuous jumps relative to previous machine learning or A.I. performance trendlines in image recognition and speech recognition; it doesn't in strategy games or natural language processing, and machine translation and arcade games are ambiguous (machine translation because metrics differ; arcade games because there is no pre-deep-learning comparison.)

Sarah Constantin, *Performance Trends In A.I.*, OTIUM (2017), <https://srconstantin.wordpress.com/2017/01/28/performance-trends-in-ai> (last visited Dec. 27, 2022). The AlphaGo Zero program discussed later, however, uses another form of deep learning—deep reinforcement learning—to excel at Go, paradigmatically thought of as a strategy game, and both translation and natural language processing have recently advanced by leaps and bounds, so this statement may no longer be true.

17. Cade Metz, "The Godfather of A.I." Leaves Google and Warns of Danger Ahead, N.Y. TIMES (May 1, 2023), <https://www.nytimes.com/2023/05/01/technology/ai-google-chatbot-engineer-quits-hinton.html>.

18. Anthony Cuthbertson, DeepMind Boss Says Human-Level AI Is Just a Few Years Away, INDEPENDENT (May 4, 2023), <https://www.independent.co.uk/tech/ai-deepmind-artificial-general-intelligence-b2332322.html>.

19. TYLER COWEN, THE GREAT STAGNATION: HOW AMERICA ATE ALL THE LOW-HANGING FRUIT OF MODERN HISTORY, GOT SICK, AND WILL (EVENTUALLY) FEEL BETTER (2011).

20. ROBERT GORDON, THE RISE AND FALL OF AMERICAN GROWTH: THE U.S. STANDARD OF LIVING SINCE THE CIVIL WAR (2016).

21. It should be noted that Cowen also makes the synchronicity point and argues more for modesty than for pessimism. He says that he and other like-minded thinkers are "relatively optimistic about the technological future of the United States, but we, along with most informed participants in these debates, are skeptical about our ability to forecast rates of economic and productivity growth many years into the future or, for that matter, even a few years ahead." So perhaps it is more accurate to call him an agnostic rather than a pessimist. Tyler Cowen, *Is Innovation Over: The Case against Pessimism*, FOREIGN AFFAIRS (Mar./Apr. 2016), <https://www.foreignaffairs.com/reviews/review-essay/2016-02-15/innovation-over>.

22. Cade Metz, Paul Allen Wants to Teach Machines Common Sense, N.Y. TIMES (Feb. 28, 2018), <https://www.nytimes.com/2018/02/28/technology/paul-allen-ai-common-sense.html>.

23. Katja Grace et al., Viewpoint: When Will AI Exceed Human Performance? Evidence from AI Experts, 62 J. A.I. RSCH. 729 (2018).

24. *Id.* at 731.

25. Mathew Barnett, Date of Artificial General Intelligence, METACULUS (Aug. 23, 2020), <https://www.metaculus.com/questions/5121/date-of-artificial-general-intelligence>.

26. METACULUS, <https://www.metaculus.com/prediction/10842/a-public-prediction-by-bryan-caplan/> <https://www.metaculus.com/prediction/10842/a-public-prediction-by-bryan-caplan/> (last visited July 13, 2023).
27. Rodney Brooks, *The Seven Deadly Sins of Predicting the Future of AI*, RODNEY BROOKS: ROBOTS, AI, & OTHER STUFF (Sep. 7, 2017), <https://rodneybrooks.com/the-seven-deadly-sins-of-predicting-the-future-of-ai>.
28. See IEEE SPECTRUM, *Human Level AI Is Right around the Corner or Hundreds of Years Away* (May 31, 2017), <https://spectrum.ieee.org/computing/software/humanlevel-ai-is-right-around-the-corner-or-hundreds-of-years-away>.
29. Rodney Brooks, *I, Rodney Brooks, Am a Robot*, 45 IEEE SPECTRUM 71, 72 (2008).
30. *Id.* at 71.
31. Vinge, *supra* note 10, at 89.
32. Mark Fischetti, *Computers versus Brains*, SCI. AM. (Nov. 1, 2011), <https://www.scientificamerican.com/article/computers-vs-brains>.
33. IEEE SPECTRUM, *supra* note 28.
34. *Id.*
35. Ajeya Cotra, *Forecasting TAI with Biological Anchors* (July 2020), <https://docs.google.com/document/d/1IJ6Sr-gPeXdSJugFulwIpvavc0atjHGM82QjfUSBGQ>.
36. Ajeya Cotra, *Two Year Update on My Personal AI Timelines*, LESSWRONG (Aug. 2, 2022), <https://www.lesswrong.com/posts/AfH2oPHCapdKicM4m/two-year-update-on-my-personal-ai-timelines>.
37. AlphaGo, DEEPMIND, <https://www.deepmind.com/research/highlighted-research/alphago> (last visited July 13, 2023).
38. David Silver et al., *Mastering the Game of Go without Human Knowledge*, 550 NATURE 354, 354 (2017).
39. *AlphaGo Zero: Starting from Scratch*, DEEPMIND (Oct. 18, 2017), <https://www.deepmind.com/blog/alphago-zero-starting-from-scratch> (emphasis added).
40. IEEE SPECTRUM, *supra* note 28.
41. Alan Turing, *Computing Machinery and Intelligence*, 49 MIND 433, 433 (1950).
42. Tyler Cowen and Michelle Dawson have argued that Alan Turing himself might not have passed the Turing Test and that the entire article is in part a meditation on the *dangers* of using imitation as our criteria. See Tyler Cowen & Michelle Dawson, *What Does the Turing Test Really Mean? And How Many Human Beings (Including Turing) Could Pass?* (June 3, 2009), <https://d101vc9winf8ln.cloudfront.net/documents/28495/original/turingfinal.pdf>. True, Turing was persecuted for being gay; bigotry can impose social distance on anyone. Cowen also theorizes that he may have had what was formerly called Asperger's Syndrome. Regardless of whether he did, "neurotypicality" would be a repellent metric to use on humans as a test of humanity. "On the spectrum" does not mean "outside the line." So, this is a nice thought experiment on the moral hazards of mimesis. Having said that, I personally found the tone and style of Turing's famous article to contradict the idea that he would have failed the

Turing Test. This is particularly true of the understated dry humor that he deploys, humor aimed quite brilliantly at anticipating emotional objections and disarming the reader of the article. That same intuitive understanding of likely skepticism, and the ability humorously to refute it, would presumably work not just on a reader but on a Turing tester, at least one with a British sense of irony.

43. John Searle, *Minds, Brains, and Programs*, 3 BEHAV. BRAIN SCIS. 417 (1980).
44. Kate Torgovnick May, *4 Talks on a Strange Phenomenon We All Experience: Consciousness*, TEDBLOG (July 22, 2013), <https://blog.ted.com/4-talks-on-a-strange-phenomenon-we-all-experience-consciousness>.
45. Alan Turing, *Computing Machinery and Intelligence*, 49 MIND 433, 447 (1950). Turing might have been surprised to find out that B. F. Skinner and the behaviorists were willing to embrace the position that humans are automata and that consciousness is an illusion and turn it into an intellectual franchise. Here is J. B. Watson, one of the founding behaviorists, writing about consciousness and the mind:

If the behaviorists are right in their contention that there is no observable mind-body problem and no observable separate entity called mind—then there can be no such thing as consciousness and its subdivision. Freud's concept borrowed from somatic pathology breaks down. There can be no festering spot in the substratum of the mind—in the unconscious—because there is no mind.

JOHN B. WATSON, THE WAYS OF BEHAVIORISM 96 (1928).
46. B.F. SKINNER, CONTINGENCIES OF REINFORCEMENT 260 (Copley Publ'g Grp. 2013) (1969).
47. It turns out to be basically true, though there is no record of the exact words used. Keith Tomson, *Huxley Wilberforce and the Oxford Museum*, AM. SCIENTIST, <https://www.americanscientist.org/article/huxley-wilberforce-and-the-oxford-museum> (last visited Feb. 17, 2024).
48. Daniel C. Dennett et al., *The Practical Requirements for Making a Conscious Robot [and Discussion]*, 349 PHIL. TRANSACTIONS: PHYSICAL SCIS. ENG'G 133, 133–36 (2023).
49. *Id.* at 135.
50. John Searle, *Our Shared Condition—Consciousness*, TEDxCERN, at 02:38 (May, 2013), https://www.ted.com/talks/john_searle_our_shared_condition_consciousness/transcript.
51. *Interview with John Searle*, NEW PHILOSOPHER (Jan. 25, 2014), <https://www.newphilosopher.com/articles/john-searle-it-upsets-me-when-i-read-the-nonsense-written-by-my-contemporaries>.
52. The italics come only from my own delighted imagination of his delivery. Various versions of the story are given. See, e.g., James Ryerson, *Sidewalk Socrates*, N.Y. TIMES (Dec. 26, 2004), <https://www.nytimes.com/2004/12/26/magazine/sidewalk-socrates.html>; Gerd Gigerenzer, *ADAPTIVE THINKING: RATIONALITY IN THE REAL WORLD* 230 (2000).
53. John R. Searle, *The Myth of the Computer*, N.Y. REV. BOOKS (Apr. 29, 1982), <https://www.nybooks.com/articles/1982/04/29/the-myth-of-the-computer>.
54. Stuart R. Hameroff, *The Brain Is Both Neurocomputer and Quantum Computer*, 31 COGNITIVE SCI. 1035 (2007); see also E. Alfinito & G. Vitiello, *The Dissipative Quantum*

Model of Brain: How Does Memory Localize in Correlated Neuronal Domains, 128 INFO. SCIS. 217 (2000); Stuart R. Hameroff, *Quantum Computation in Brain Microtubules? The Penrose-Hameroff "Orch OR" Model of Consciousness*, 356 PHIL. TRANSACTIONS ROYAL SOC'Y LONDON 1869 (1998); ROGER PENROSE, THE EMPEROR'S NEW MIND: CONCERNING COMPUTERS, MINDS AND THE LAWS OF PHYSICS (1989).

55. A. Litt et al., *Is the Brain a Quantum Computer?*, 30 COGNITIVE SCI. 593, 593 (2006).

Scientific attempts to understand human thinking have historically drawn on analogies with contemporary technologies, from clockworks to telephone switchboards to digital computers. Today, one of the most exciting emerging technologies is quantum computation, which attempts to overcome limitations of classical computers by employing phenomena unique to quantum-level events, such as nonlocal entanglement and superposition. It is therefore not surprising that many researchers have conjectured that quantum effects in the brain are crucial for explaining psychological phenomena, including consciousness. We argue, however, that explaining brain function by appeal to quantum mechanics is akin to explaining bird flight by appeal to atomic bonding characteristics. . . . We conclude that understanding brain function is unlikely to require quantum computation or similar mechanisms. *Id.* at 593–94.

56. For an introduction to illusionism, and an attempt to rebut its principal detractors, see Keith Frankish, *Illusionism as a Theory of Consciousness*, 23 J. CONSCIOUSNESS STUDS. 11 (2016), <https://philpapers.org/rec/FRAIAA-4>.

57. Michael Rescorla, *The Computational Theory of Mind*, STAN. ENCYC. PHIL. (Edward N. Zalta ed., 2020), <https://plato.stanford.edu/entries/computational-mind>.

58. Giulio Tononi, *An Information Integration Theory of Consciousness*, 5 BMC NEUROSCIENCE, no. 42, 2004, <https://bmcbiochem.biomedcentral.com/articles/10.1186/1471-2202-5-42>.

59. *Id.*

60. Mariana Lenharo, *Prominent Consciousness Theory Is Slammed as Bogus Science*, SCI. AM. NATURE MAG. (Sept. 21, 2023), <https://www.scientificamerican.com/article/prominent-consciousness-theory-is-slammed-as-bogus-science>.

61. Cogitate Consortium et al., An Adversarial Collaboration to Critically Evaluate Theories of Consciousness (June 26, 2023) (unpublished manuscript), <https://www.biorxiv.org/content/10.1101/2023.06.23.546249v1.full.pdf>. An entertaining account of the test, and its accompanying wager in wine, can be found at Elizabeth Finkel, *"Adversarial" Search for Neural Basis of Consciousness Yields First Results*, SCI. (June 25, 2023), <https://www.science.org/content/article/search-neural-basis-consciousness-yields-first-results>.

62. Finkel, *supra* note 61.

63. Tononi, *supra* note 58 (emphasis added).

64. PATRICK BUTLIN & ROBERT LONG ET AL., CONSCIOUSNESS IN ARTIFICIAL INTELLIGENCE: INSIGHTS FROM THE SCIENCE OF CONSCIOUSNESS 13 (2023), <https://arxiv.org/pdf/2308.08708.pdf>.

65. Computational functionalism's roots can be traced back to the work of Hilary Putnam, whose opinions about the theory changed over time. See Oron Shagrir, *The Rise and Fall of Computational Functionalism*, in HILARY PUTNAM 220 (Yemima

- Ben-Menahem ed., Cambridge Univ. Press 2005), <https://www.cambridge.org/core/books/abs/hilary-putnam/rise-and-fall-of-computational-functionalism>.
66. Oliver Whang, *How to Tell if Your A.I. Is Conscious*, N.Y. TIMES (Sept. 18, 2023), <https://www.nytimes.com/2023/09/18/science/ai-computers-consciousness.html>.
67. Matthias Michel & Hakwan Lau, Commentary, *Higher-Order Theories Do Just Fine*, 12 COGNITIVE NEUROSCIENCE 77, 78 (2021), <https://www.tandfonline.com/doi/abs/10.1080/17588928.2020.1839402>.
68. BUTLIN & LONG ET AL., *supra* note 64.
69. Whang, *supra* note 66.
70. Corporations are a partial exception, as I will explain later.
71. Irving John Good, *Speculations on the First Ultraintelligent Machine*, 6 ADVANCES COMPUTS. 31 (1966).
72. Stephen Hawking et al., *Stephen Hawking: Transcendence Looks at the Implications of Artificial Intelligence: But Are We Taking AI Seriously Enough?*, INDEPENDENT (May 1, 2014), <https://www.independent.co.uk/news/science/stephen-hawking-transcendence-looks-at-the-implications-of-artificial-intelligence-but-are-we-taking-ai-seriously-enough-9313474.html>
73. Vinge, *supra* note 10, at 92.
74. NICK BOSTROM, SUPERINTELLIGENCE: PATHS, DANGERS, STRATEGIES (2014).
75. Catherine Clifford, *Elon Musk: “Mark My Words—A.I. Is Far More Dangerous Than Nukes,”* CNBC (Mar. 14, 2018), <https://www.cnbc.com/2018/03/13/elon-musk-at-sxsw-a-i-is-more-dangerous-than-nuclear-weapons.html>.
76. Cade Metz, *Mark Zuckerberg, Elon Musk and the Feud over Killer Robots*, N.Y. TIMES (June 9, 2018), <https://www.nytimes.com/2018/06/09/technology/elon-musk-mark-zuckerberg-artificial-intelligence.html>.
77. Chris Williams, *AI Guru Ng: Fearing a Rise of Killer Robots Is Like Worrying about Overpopulation on Mars*, REGISTER (Mar. 19, 2015), https://www.theregister.com/2015/03/19/andrew_ng_baidu_ai.
78. James Vincent, *Elon Musk and Top AI Researchers Call for Pause on “Giant AI Experiments,”* VERGE (Mar. 29, 2023), <https://www.theverge.com/2023/3/29/23661374/elon-musk-ai-researchers-pause-research-open-letter>.
79. *Statement on AI Risk*, CTR. FOR AI SAFETY, <https://www.safe.ai/statement-on-ai-risk> (last visited July 12, 2023).
80. James Vincent, *Top AI Researchers and CEOs Warn against “Risk of Extinction” in 22-Word Statement*, VERGE (May 30, 2023), <https://www.theverge.com/2023/5/30/23742005/ai-risk-warning-22-word-statement-google-deepmind-openai>.
81. BOSTROM, *supra* note 74, at i.
82. *Id.* at 128.
83. BOSTROM, *supra* note 74, at 123.

84. *Id.* at 134. The original idea came from Marvin Minsky. STUART RUSSELL & PETER NORVIG, ARTIFICIAL INTELLIGENCE: A MODERN APPROACH 1039 (3d ed. 2010).
85. See Eliezer Yudkowsky, *Complex Value Systems in Friendly AI*, in ARTIFICIAL GENERAL INTELLIGENCE 388 (2011).
86. WIKIPEDIA, s.v. *Snail Darter*, https://en.wikipedia.org/wiki/Snail_darter (last visited July 12, 2023).
87. Tennessee Valley Authority v. Hiram Hill et al., 437 U.S. 153 (1978).
88. I owe the phrase to Joel Shepherd's novels.
89. Eliezer Yudkowsky, *Pausing AI Developments Isn't Enough. We Need to Shut It All Down*, TIME MAG. (March 29, 2023), <https://time.com/6266923/ai-eliezer-yudkowsky-open-letter-not-enough/>.
90. Yudkowsky, *supra* note 11.
91. RUSSELL & NORVIG, *supra* note 84, at 3.
92. See Cowen & Dawson, *supra* note 42.
93. See Kevin Roose, *An A.I.-Generated Picture Won an Art Prize. Artists Aren't Happy*, N.Y. TIMES: THE SHIFT (Sep. 2, 2022), <https://www.nytimes.com/2022/09/02/technology/ai-artificial-intelligence-artists.html> (reporting that an artwork generated by AI won first place in the emerging digital artists' contest at the Colorado State Fair's annual art competition). Not everyone has joined the chorus of lamentation. Farhad Manjoo argues persuasively that AI may actually benefit both art and artists. "What accounts for my sunny stance? History offers one clue: Technologies that made art easier to produce have rarely ended up stifling human creativity. Electronic synthesizers didn't eliminate the need for people who play musical instruments. Auto-Tune didn't make singing on pitch obsolete. Photography didn't kill painting, and its digitization didn't obviate the need for professional photographers." Farhad Manjoo, *A Creator (Me) Made a Masterpiece With A.I.*, N.Y. TIMES (Aug 25, 2023), <https://www.nytimes.com/2023/08/25/opinion/ai-art-intellectual-property.html>. He makes a similar argument about the effect of AI on computer programming—abstracting from the precise *form* of the activity as it is currently constituted to the *role* that humans play. Coding may cease to be an activity humans engage in, but humans will still have a role in working to instruct computers on how to perform complex tasks, even if it is increasingly done through natural language. Farhad Manjoo, *It's the End of Computer Programming as We Know It. (And I Feel Fine)*, N.Y. TIMES (June 2, 2023), <https://www.nytimes.com/2023/06/02/opinion/ai-coding.html>.
94. This response locates the creativity in the human being who forms the prompts. We do not doubt Ansel Adams's artistry even though it was conveyed through a camera rather than a paintbrush. Why should this be any different? The difficulty is the scale of the relative contributions of human and machine. While some prompts exhaustively detail style, composition, shadow, exposure, and even type of camera lens, others consist of only two or three words. The US Copyright Office attempts to navigate just this line: "In the Office's view, it is well-established that copyright can protect only material that is the product of human creativity. Most fundamentally,

the term ‘author,’ which is used in both the Constitution and the Copyright Act, excludes non-humans.” 37 C.F.R. § 202 (2023). Thus, the Copyright Office will register copyrights over AI-generated material only if the human being contributes significant artistic input. For example, a comic book that contained human text and AI-generated images was copyrightable as a whole, but the individual machine-generated images in the comic were not. *See id.* A similar approach has been taken in patent law, with courts declaring that only a human can be an inventor. *Thaler v. Vidal* 43 F.4th 1207 (Fed. Cir. 2022).

95. Roland Barthes, *The Death of the Author*, in IMAGE, MUSIC, TEXT (Stephen Heath trans., 1977), <https://sites.tufts.edu/english292b/files/2012/01/Barthes-The-Death-of-the-Author.pdf>.

96. Baumol’s cost disease or cost effect is a description of a counterintuitive economic phenomenon. Some occupations or economic fields experience very little productivity growth over time, yet potentially share in the higher wages that more productive occupations receive. You still need four musicians and 90 minutes for your eighteenth-century string quartet, just as you did when it was written. Meanwhile average labor productivity has dramatically increased. WIKIPEDIA, s.v. *Baumol Effect*, https://en.wikipedia.org/wiki/Baumol_effect (last visited July 12, 2023).

97. Academic scholarship on the issue is split between will theorists, who believe the predicate for rights is the ability freely to make rational moral choices, and interest theorists, who believe that rights should be given to all of those with the relevant moral interests, even if, like a baby or a person deemed legally insane, they have no such ability. I identify more with the interest theory side of the debate, but my point here is only that the criteria we apply to AI will reflect our underlying moral presuppositions and that any test for personhood will probably include strands reflecting both lines of thought. Our folkways are less methodologically monocultural than our academic theories and, from my point of view, that is not always a bad thing. For further discussion, see MATTHEW H. KRAMER ET AL., A DEBATE OVER RIGHTS (1998); DAVID GUNKEL, ROBOT RIGHTS (2018).

98. The AI doomers, who already think that mindless implementation of human directives might doom the species, would view the valorization of both autonomy and machine societies as lunacy on stilts.

99. *See, e.g.*, Alan F. T. Winfield, *When Robots Tell Each Other Stories: The Emergence of Artificial Fiction*, in NARRATING COMPLEXITY 39 (R. Walsh & S. Stepney eds., 2017), https://link.springer.com/chapter/10.1007/978-3-319-64714-2_4.

100. Sourya Acharya & Samarth Shukla, *Mirror Neurons: Enigma of the Metaphysical Modular Brain*, 3 J. NAT'L. SCI. BIOLOGY & MED. 118, 119 (2012).

101. *See, e.g.*, Noam Chomsky, *On Cognitive Structures and Their Development: A Reply to Piaget*, in LANGUAGE AND LEARNING: THE DEBATE BETWEEN JEAN PIAGET AND NOAM CHOMSKY (Massimo Piattelli-Palmarini ed., 1980); SHAUN GALLAGHER, HOW THE BODY SHAPES THE MIND (2005); DAN ZAHAVI, SUBJECTIVITY AND SELFHOD: INVESTIGATING THE FIRST-PERSON PERSPECTIVE (2005); EVAN THOMPSON, MIND IN LIFE: BIOLOGY, PHENOMENOLOGY, AND THE SCIENCES OF MIND (2010). For an accessible overview of the broad and interdisciplinary work

in “embodied cognition,” see Lawrence Shapiro & Shannon Spaulding, *Embodied Cognition*, STAN. ENCYC. PHIL. (Edward N. Zalta ed., 2021), <https://plato.stanford.edu/archives/win2021/entries/embodied-cognition>.

102. GEORGE LAKOFF & MARK JOHNSON, PHILOSOPHY IN THE FLESH: THE EMBODIED MIND AND ITS CHALLENGE TO WESTERN THOUGHT (1999).

103. Harry Lambert, *Is AI a Danger to Humanity or Our Salvation?*, NEW STATESMAN (June 21, 2023), <https://www.newstatesman.com/long-reads/2023/06/men-made-future-godfathers-ai-geoffrey-hinton-yann-lecun-yoshua-bengio-artificial-intelligence>. I am grateful to Tyler Cowen and Alex Tabarrok’s site *Marginal Revolution* for highlighting this article and many other useful pieces of information. MARGINAL REVOLUTION, <https://marginalrevolution.com> (last visited June 13, 2023).

104. See LUDWIG WITTGENSTEIN, PHILOSOPHICAL INVESTIGATIONS para. 280, at 96e (G. E. M. Anscombe, 1958), https://archive.org/details/philosophicalinvestigations_201911/page/n101:

Someone paints a picture in order to shew how he imagines a theatre scene. And now I say: “This picture has a double function: it informs others, as pictures or words inform—but for the one who gives the information it is a representation (or piece of information?) of another kind: for him it is the picture of his image, as it can’t be for anyone else. To him his private impression of the picture means what he has imagined, in a sense in which the picture cannot mean this to others.”—And what right have I to speak in this second case of a representation or piece of information—if these words were rightly used in the first case?

See also LUDWIG WITTGENSTEIN, LAST WRITINGS ON THE PHILOSOPHY OF PSYCHOLOGY: THE INNER AND THE OUTER §84 (G.H. von Wright et al. eds., vol. 2, 1992) (“[t]he ‘inner’ is a delusion. That is: the whole complex of ideas alluded to by this word is like a painted curtain drawn in front of the scene of the actual word use”). Some of Wittgenstein’s discussion seems to me to imply that there is indeed a rich inner life, but that it is one to which language and behavior—the “painted curtain”—can give others only limited and uncertain access. “There’s no art to find the mind’s construction in the face.” WILLIAM SHAKESPEARE, MACBETH Act 4, sc. 1. I have no clear view of your inner theater, merely of the reduced and possibly falsified, distorted, or misunderstood version of it presented by the “painted curtain” used in traditional stage productions, dramatizing certain aspects of the play. In that sense, Wittgenstein is fulfilling his characteristic role: showing us how misplaced terminology and reliance on reified definitions can cause us to understand something poorly, if at all. For an excellent discussion of these points, see Paul Standish, *Inner and Outer, Psychology and Wittgenstein’s Painted Curtain*, 56 J. PHIL. EDUC. 115 (2022). I would happily acknowledge, however, that one reason Wittgenstein is of such enduring interest is that his words can be used as “intuition pumps” for a wide range of ideas, not all of them consistent.

105. Neither theory has prevailed, and each have been shown to explain some things well but others poorly. See, e.g., Mariana Lenharo, *Decades-Long Bet on Consciousness Ends—and It’s Philosopher 1, Neuroscientist 0*, NATURE (June 24, 2023), <https://www.nature.com/articles/d41586-023-02120-8>; Rufin VanRullen & Ryota Kanai, *Deep Learning and the Global Workspace Theory*, 44 TRENDS NEUROSCIENCES 692 (2021).

106. RAY KURZWEIL, HOW TO CREATE A MIND: THE SECRET OF HUMAN THOUGHT REVEALED 213 (2012).
107. See *Report on a Comprehensive European Industrial Policy on Artificial Intelligence and Robotics*, at 37–40 (Jan. 1, 2019), https://www.europarl.europa.eu/doceo/document/A-8-2019-0019_EN.pdf (discussing potential legal liability for AI).
108. See, e.g., Brandon Garrett & Cynthia Rudin, *The Right to a Glass Box: Rethinking the Use of Artificial Intelligence in Criminal Justice*, 109 CORNELL L. REV. (forthcoming 2023).
109. Steven Zeitchik, *Is Artificial Intelligence about to Transform the Mammogram?*, WASH. POST (Dec. 21, 2021), <https://www.washingtonpost.com/technology/2021/12/21/mammogram-artificial-intelligence-cancer-prediction>.
110. My colleagues Brandon Garrett and Cynthia Rudin argue that we have accepted all too easily the trade-off between transparency and accuracy, one they believe to be illusory. They claim that in many cases far simpler, human legible criteria can produce results that are equally accurate. Garrett & Rudin, *supra* note 108. Of course, the legal system in particular should require algorithms that are open and transparent rather than closed and proprietary and Garret and Rudin's specific criminal justice arguments are extremely persuasive. I remain unconvinced, however, of the larger claim that the trade-off will always be as illusory as they suggest. Many of the stories of machine learning's greatest successes, such as that of breast cancer detection, reinforce this conclusion.
111. Tyler Cowen, *The Taxman Will Eventually Come for AI, Too*, BLOOMBERG (Apr. 17, 2023), <https://www.bloomberg.com/opinion/articles/2023-04-17/the-taxman-will-eventually-come-for-ai-too>.
112. Lawrence B. Solum, *Legal Personhood for Artificial Intelligences*, 70 N.C. L. REV. 1231, 1286–87 (1992).

CHAPTER 3

1. Felix Cohen, *Transcendental Nonsense and the Functional Approach*, 35 COLUM. L. REV. 809, 813 (1935).
2. Trustees of Dartmouth College v. Woodward, 17 U.S. 518, 636 (1819).
3. See, e.g., TYLER COWEN, BIG BUSINESS: A LOVE LETTER TO AN AMERICAN ANTIHERO (2019).
4. *Draft Report with Recommendations to the Commission on Civil Law Rules on Robotics* 12 (May 31, 2016), https://www.europarl.europa.eu/doceo/document/JURI-PR-582443_EN.pdf.
5. James Vincent, *Giving Robots “Personhood” Is Actually about Making Corporations Accountable*, VERGE (Jan. 19, 2017), <https://www.theverge.com/2017/1/19/14322334/robot-electronic-persons-eu-report-liability-civil-suits>.
6. *Dartmouth College*, 17 U.S. at 636.
7. Citizens United v. FEC, 558 U.S. 310 (2010).
8. Teneille R. Brown, *In-Corp-O-Real: A Psychological Critique of Corporate Personhood and Citizens United*, 12 FLA. ST. UNIV. BUS. REV. 1, 1 (2013).

9. The Case of Sutton's Hosp. (1612) 77 Eng. Rep. 960, 970 (K.B.).
10. Katsuhito Iwai, *Persons, Things and Corporations: The Corporate Personality Controversy and Comparative Corporate Governance*, 47 AM. J. COMPAR. L., 583, 583–84 (1999).
11. A.V. Dicey, *The Combination of Laws as Illustrating the Relation between Law and Opinion in England During the Nineteenth Century*, 17 HARV. L. REV. 511, 513 (1904).
12. Morton J. Horwitz, *Santa Clara Revisited: The Development of Corporate Theory*, 88 W. VA. L. REV. 173, 180 (1986).
13. Cohen, *supra* note 1, at 811.
14. John Dewey, *The Historic Background of Corporate Legal Personality*, 35 YALE L.J. 655, 656 (1926).
15. THOMAS HOBBS, THE LEVIATHAN 31 (G.A.J. Rogers & Karl Schuhmann eds., Bloomsbury Acad. 2006) (1651). One could take Hobbes's argument further. Is money "the money of fools"? Money, too, relies on a socially constructed agreement that these grubby pieces of drab green paper contain exchangeable value. When *that* social agreement breaks down, as in Weimar Republic inflation, the consensually created value is suddenly thrown into radical doubt.
16. See Phillip Blumberg, *The Corporate Entity in an Era of Multinational Corporations*, 15 DEL. J. CORP. L. 283, 293–95 (1990) (discussing the emergence of the associational or aggregate theory of corporate personality in American jurisprudence); see also Michael J. Phillips, *Reappraising the Real Entity Theory of the Corporation*, 21 FLA. ST. UNIV. L. REV. 1061, 1065–67 (1994) (discussing the history and development of the associational or aggregate theory of corporate personality).
17. Horwitz, *supra* note 12, at 178; (1986) (quoting Argument for Defendant, *San Mateo v. Southern Pacific Railroad Co.*, 116 U.S. 138 (1882) (collected in *Cases and Points* at 10 [available in Harvard Law School Library]) (emphasis omitted)).
18. M.C. Jensen & W.H. Meckling, *Agency Costs and the Theory of the Firm*, in CORPORATE GOVERNANCE: VALUES, ETHICS AND LEADERSHIP 83 (Robert Ian Tricker ed., 2019) (emphasis added).
19. *Id.*
20. Often attributed, without much evidence, to Samuel Clemens, a.k.a. Mark Twain.
21. In practice, the legal status of slaves was more complex, but the basic point can hardly be denied. See generally Mark Tushnet, *The American Law and Slavery, 1810–1860: A Study in the Persistence of Legal Autonomy*, 10 L. & SOCY REV. 119 (1975) (discussing the limited legal protections and rights that existed for slaves); see also Visa Kurki, *Animals, Slaves, and Corporations: Analyzing Legal Thinghood*, 18 GER. L.J. 1069, 1086 ("slaves in the antebellum United States held certain rights toward both their owners and third parties"); THOMAS D. MORRIS, SOUTHERN SLAVERY AND THE LAW, 1619–1860 (1996) (documenting the tensions between the common law and slavery).
22. *Meet the Press, McCain: Citizens United Decision an "Outrage"* (NBC television broadcast Jan. 28, 2012), <https://www.nbcnews.com/video/mccain-citizens-united-decision-an-outrage-43863107634>.

23. *Citizens United v. FEC*, 558 U.S. 310, 466 (2010) (Stevens, J., dissenting).
24. See, e.g., RALPH C. FERRARA ET AL., SHAREHOLDER DERIVATIVE LITIGATION: BESIEGING THE BOARD 10–54 (2013); John C. Coffee, Jr., “*No Soul to Damn: No Body to Kick*”: An *Unscandalized Inquiry into the Problem of Corporate Punishment*, 79 MICH. L. REV. 386, 386 (1981).
25. See, e.g., JOHN POYNDER, LITERARY EXTRACTS 268 (1844), <https://archive.org/details/literaryextracts01poynuoft/page/268>.
26. *First Nat'l Bank of Boston v. Bellotti*, 435 U.S. 765, 822–26 (1978) (Rehnquist, J., dissenting) (emphasis added).
27. *Id.* at 784.
28. *Id.* at 802 (Burger, J., concurring). Justice Scalia's concurrence in *Citizens United* made a similar point, using a version of association theory to argue that the freedom of speech protected by the First Amendment includes the “freedom to speak in association with other individuals, including association in the corporate form.” I am grateful to Larry Lessig for stressing the importance of this point.
29. There are a number of wonderful historical accounts of the history of the struggle for constitutional personhood for corporations, including the roles of Conkling and of the Supreme Court’s reporter. The story begins with the work of Charles and Mary Beard, often presented—perhaps unfairly—as the proponents of the “conspiracy theory” of the Fourteenth Amendment. The most definitive historical account of Conkling’s argument and of the actual framing of the Amendment is HOWARD J. GRAHAM, EVERYMAN’S CONSTITUTION: HISTORICAL ESSAYS ON THE FOURTEENTH AMENDMENT, THE “CONSPIRACY THEORY,” AND AMERICAN CONSTITUTIONALISM (1968). C. PETER MAGRATH, MORRISON R. WAITE: THE TRIUMPH OF CHARACTER (1963) should be credited as pointing me toward invaluable primary sources and providing important context. THOM HARTMANN, “UNEQUAL PROTECTION”: HOW CORPORATIONS BECAME “PEOPLE”—AND HOW YOU CAN FIGHT BACK (2010) offers a passionate negative case against corporate constitutional rights, studded with a surprising amount of historical detail and facsimiles of the original documents. Last, but by no means least, I am indebted to Adam Winkler’s extremely readable WE, THE CORPORATIONS: HOW AMERICAN BUSINESSES WON THEIR CIVIL RIGHTS (2018). I disagree with some of Winkler’s conclusions and emphases but wholeheartedly recommend the book.
30. Guy Gugliotta, *New Estimate Raises Civil War Death Toll*, N.Y. TIMES (Apr. 2, 2012), <https://www.nytimes.com/2012/04/03/science/civil-war-toll-up-by-20-percent-in-new-estimate.html>.
31. U.S. CONST. amend. XIV, § 1.
32. Slaughter-House Cases, 83 U.S. 36, 37 (1872).
33. Ralph Waldo Emerson, *The Problem*, POETS.ORG, <https://poets.org/poem/problem> (last visited Sept. 22, 2023).
34. Charles Beard, *Corporations and Natural Rights*, 12 VA. Q. REV. 337 (1936), reprinted in CHARLES BEARD, JEFFERSON, CORPORATIONS AND THE CONSTITUTION 1, 23 (1936).

35. HOWARD J. GRAHAM, "*Buidled Better Than They Knew*": *The Framers, the Railroads and the Fourteenth Amendment*, in EVERYMAN'S CONSTITUTION 443 (Wis. Hist. Soc'y 2013) (1968). The bracketed comment is my own, but I believe it fairly represents Graham's argument.
36. *Id.* at 394.
37. HARTMANN, *supra* note 29, at 47–48. Winkler cites Graham for the same proposition, but I could find no reference to this fact on the cited page. WINKLER, *supra* note 29, at 152.
38. JACK BEATTY, AGE OF BETRAYAL: THE TRIUMPH OF MONEY IN AMERICA, 1865–1900 172 (2007).
39. See e.g., WINKLER, *supra* note 29, at 152.
40. *Id.* at 157–58.
41. CHARLES WALLACE COLLINS, THE FOURTEENTH AMENDMENT AND THE STATES 41 (1912).
42. See, e.g., WINKLER, *supra* note 29.
43. Jesse J. Holland, *Scalia Says No "Falling Out" with Roberts*, YAHOO! NEWS (July 20, 2012), <https://news.yahoo.com/scalia-says-no-falling-roberts-005026260.html>.
44. HOWARD J. GRAHAM, *The "Conspiracy Theory" of the Fourteenth Amendment: Part II*, in EVERYMAN'S CONSTITUTION 93–94 (Wis. Hist. Soc'y 2013) (1968) (emphasis added).
45. First Nat'l Bank of Boston v. Bellotti, 435 U.S. 765, 826 (1978) (Rehnquist, J., dissenting).
46. Nat'l Mut. Ins. Co. of D.C. v. Tidewater Transfer Co., 337 U.S. 582, 646 (1949) (Frankfurter, J., dissenting).
47. See Wheeling Steel Corp. v. Glander, 337 U.S. 562, 576–81 (1949) (Douglas, J., dissenting).
48. Citizens United v. FEC, 558 U.S. 310 (2010).
49. There are lots of other reasons to doubt originalism. See, e.g., James Boyle, *A Process of Denial: Bork and Post-Modern Conservatism*, 3 YALE L. & HUMANS. 263 (1991); Mitchell Berman, *Originalism Is Bunk*, 84 N.Y. L. REV. 1 (2009).
50. Obergefell v. Hodges, 576 U.S. 644, 647 (2015).
51. *Id.* at 673.
52. *Citizens United*, 558 U.S. at 310. There are a plethora of critics of the Court's decision in *Citizens United*. See, e.g., Timothy K. Kuhner, *The Democracy to Which We Are Entitled: Human Rights and the Problem of Money in Politics*, 26 HARV. HUM. RTS. J. 39, 43–44 (2013) (discussing how campaign spending can cause corruption in democracies); Ciara Torres-Spelliscy, *Safeguarding Markets from Pernicious Pay to Play: A Model Explaining Why the SEC Regulates Money in Politics*, 12 CONN. PUB. INT. L.J. 361, 362–63 (2013) (arguing that the decision caused harm to investors and securities markets); Monica Youn, *Small-Donor Public Financing in the Post-Citizens United Era*, 44 J. MARSHALL L. REV. 619, 620 (2011) (criticizing an accountability crisis due to the decision); Matthew A. Melone, *Citizens United and Corporate Political Speech: Did the Supreme Court Enhance Political Discourse or Invite Corruption?*, 60 DEPAUL L. REV. 29, 88 (2010) (criticizing the Court's reasoning as unpersuasive).

CHAPTER 4

1. *Client: Tommy*, NONHUMAN RIGHTS GROUP, <https://www.nonhumanrights.org/client-tommy> (last visited Jan. 1, 2023).
2. *Who We Are*, NONHUMAN RIGHTS GROUP, <https://www.nonhumanrights.org/who-we-are> (last visited Jan. 1, 2023).
3. I profited greatly from Kate Darling's excellent recent book, *The New Breed*, which explicitly makes this analogy. KATE DARLING, THE NEW BREED: WHAT OUR HISTORY WITH ANIMALS REVEALS ABOUT OUR FUTURE WITH ROBOTS (2021).
4. Many philosophers have used consciousness to distinguish humans from other creatures. See, e.g., Immanuel Kant, *Anthropology from a Pragmatic Point of View*, in CAMBRIDGE EDITION OF THE WORKS OF IMMANUEL KANT 227, 239 (Robert Louden ed., Gunter Zoller trans., 2010) (1798) (arguing that humans' ability to have the representation of "I" raises them above other beings); Daniel Dennet, *Animal Consciousness: What Matters and Why*, 62 SOC. RSCH. 691, 702–03 (1995) (arguing that human consciousness is a kind of "user-illusion" that entails an "informational organization" not present in other species, although he pleads mischaracterization when being described as arguing that "other species lack our kind of *self-consciousness*"). However, other philosophers focus on developing criteria for determining personhood rather than on distinguishing humans from other beings. See, e.g., David DeGrazia, *Human-Animal Chimeras: Human Dignity, Moral Status, and Species Prejudice*, 38 METAPHILOSOPHY 309, 319–20 (2007) (arguing that personhood is not coextensive with humanity but that the being needs a "sufficiently complex" form of consciousness that has a "high enough" degree of autonomy, rationality, self-awareness, linguistic competence, sociability, moral agency, and intentionality in action). Other philosophers reject personhood as the primary moral criteria. See, e.g., JEREMY BENTHAM, AN INTRODUCTION TO THE PRINCIPLES OF MORALS AND LEGISLATION, ch. xvii, at 351 (1789), <http://www.koeblergerhard.de/Fontes/BenthamJeremyMoralsandLegislation1789.pdf> (making an early and forceful argument for a sentience-oriented view of morality and identifying the capacity to suffer as the relevant criteria).
5. True, we might think that fairness requires that we should treat all property owners in the same way. Therefore if we allow corporations to hold title to property in the same way that natural persons do, they should be given the same due process protections. But this is a secondary consequence of a classification we know at all times is merely a convenient fiction, and it is one that we could forgo or modify without violating a central tenet of the moral duties we believe to be owed to natural persons.
6. The hard-core utilitarian might argue that all moral claims reduce to some kind of utilitarian, consequentialist analysis and thus that the *categorical* nature of the distinction I am making is illusory. In the abstract, that is a fair point. Concretely, I could debate the philosophy at length, but I think the crispest response here is "bullshit." If you honestly tell me that, beyond the level of formal deontological classification, your framework for assessing the moral rights and duties of your child is the same as that for assessing the rights and duties of IBM, there is little I can do

apart from shaking my head and calling social services. And possibly invoking the shade of Bernard Williams.

7. Petitioners' Memorandum of Law in Support of Order to Show Cause & Writ of Habeas Corpus and Order Granting the Immediate Release of Tommy at 1, People ex rel. Nonhuman Rts. Project, Inc. v. Lavery, 124 A.D.3d 149 (N.Y. App. Div. 2014) (No. 518336) [hereinafter Petitioners' Memorandum], <https://www.nonhumanrights.org/content/uploads/Memorandum-of-Law-Tommy-Case.pdf>.

8. Felix Cohen, *Transcendental Nonsense and the Functional Approach*, 35 COLUM. L. REV. 809, 813 (1935).

9. A number of scholars have made this point. The canonical reference is Darling's recent work, *supra* note 3.

10. See, e.g., BENTHAM, *supra* note 4.

11. ARISTOTLE, POLITICS bk. 1, pt. II (Benjamin Jowett trans.) (350 BCE), <http://classics.mit.edu/Aristotle/politics.mb.txt>.

12. Genesis 1:26 (New Revised Standard Version, Anglicised). The Qur'an takes a similar line: "It is God who provided for you all manner of livestock, that you may ride on some of them and from some you may derive your food. And other uses in them for you to satisfy your heart's desires. It is on them, as on ships, that you make your journeys." Qur'an 40:79–80.

13. See, e.g., Bronwyn Finnigan, *Buddhism and Animal Ethics*, 12 PHIL. COMPASS 7 e12424 (2017); GANANATH OBEYESEKERE, IMAGINING KARMA: ETHICAL TRANSFORMATION IN AMERINDIAN, BUDDHIST, AND GREEK REBIRTH (2002); Paul Waldau, *Buddhism and Animal Rights*, in CONTEMPORARY BUDDHIST ETHICS 81 (Damien Keown ed., 2013).

14. It turns out to be basically true, though it may have the identity of the original user of the metaphor incorrect. Keith Tomson, *Huxley Wilberforce and the Oxford Museum*, AM. SCIENTIST, <https://www.americanscientist.org/article/huxley-wilberforce-and-the-oxford-museum> (last visited Feb. 17, 2024).

15. CHARLES DARWIN, THE DESCENT OF MAN AND SELECTION IN RELATION TO SEX 104–05 (Barnes & Noble Books, 2004) (1871).

16. *Id.* at 105.

17. See, e.g., Irene M. Pepperberg, *Animal Language Studies: What Happened?*, 24 PSYCHONOMIC BULL. & REV. 181 (2017) (a personal history of the research into animals' communication systems); Frans de Waal, *The Surprising Complexity of Animal Memories*, ATLANTIC (June 2, 2019), <https://www.theatlantic.com/science/archive/2019/06/surprising-complexity-animal-memories/589420> ("[c]himpanzees, birds, and even rats have shown signs of reviewing their own past to prepare for the future"); FRANS DE WAAL, THE BONOBO AND THE ATHEIST (2013) (arguing that empathy and altruism predate religion and coevolved in nonhuman primates while referencing studies purporting to show benevolence in other species); Mary Bates, *Problem-Solving Parrots Understand Cause and Effect*, WIRED (Oct. 17, 2013), <https://www.wired.com/2013/10/problem-solving-parrots-understand-cause-and-effect> (summarizing a study showing that parrots can understand cause and effect); Jessica Pierce, *Do Animals Experience Grief?*,

SMITHSONIAN MAG.: SCI. (Aug. 24, 2018), <https://www.smithsonianmag.com/science-nature/do-animals-experience-grief-180970124> (discussing research that provides evidence that animals are aware of death and will sometimes mourn and ritualize their dead); Peter Mcgraw & Joel Warner, *Do Animals Have a Sense of Humour?*, NEW SCIENTIST (Mar. 27, 2014), <https://www.newscientist.com/article/dn25312-do-animals-have-a-sense-of-humour> (discussing studies indicating that nonhuman animals have senses of humor).

18. FRANS DE WAAL, ARE WE SMART ENOUGH TO KNOW HOW SMART ANIMALS ARE? 4 (2016).
19. MARK TWAIN, FOLLOWING THE EQUATOR 238 (Harper & Bros. 1925), *quoted in* DE WAAL, *supra* note 18, at 14.
20. DE WAAL, *supra* note 18.
21. *Id.* at 5.
22. SAMUEL BUTLER, EREWHON: OR, OVER THE RANGE 144 (The Project Gutenberg ed., 2005) (1872), <https://ia601002.us.archive.org/19/items/E4CS4/Erewhon.pdf>.
23. Marc Hauser, *Origin of the Mind*, 301 SCI. AM. 44, 44–45 (2009). It is worth noting that Hauser was reportedly found guilty of research misconduct by the DHSS Office of Research Integrity and Harvard's Faculty of Arts and Sciences. See Marc Hauser “Engaged in Research Misconduct,” HARV. MAG. (Sept. 5, 2012). While those findings relate to experiments I do not discuss here, they put an obvious caution flag over his arguments.
24. Hauser, *supra* note 23, at 46.
25. The phrase came to us via Hitchcock’s first “talkie,” the 1929 *Blackmail*. Apparently it was coined by New York architect Stanford White, “who purportedly used the line to induce women whom he wanted to seduce to visit the townhouse that he had furnished with etchings of nudes.” Pascal Treguer, *History of Come Up and See My Etchings*, WORD HISTS., <https://wordhistories.net/2020/01/11/come-see-my-etchings> (last visited July 18, 2023). Psychologists still use it as a classic example of indirect or strategic speech:

Speakers often do not state requests directly but employ innuendos such as “Would you like to see my etchings?” Though such indirectness seems puzzlingly inefficient, it can be explained by a theory of the strategic speaker, who seeks plausible deniability when he or she is uncertain of whether the hearer is cooperative or antagonistic.

James J. Lee & Stephen Pinker, *Rationales for Indirect Speech: The Theory of the Strategic Speaker*, 117 PSYCH. REV. 785, 785 (2010). All of which goes to show that having psychologists explain your pickup lines is an infallible contraceptive.
26. Hauser, *supra* note 23, at 51.
27. Petitioners’ Memorandum, *supra* note 7, at 1, 4.
28. See, e.g., Arthur E. Brown, *Grief in the Chimpanzee*, 13 AM. NATURALIST 173 (1879) (documenting displays of grief by a chimpanzee in captivity); Elizabeth V. Lonsdorf et al., *Why Chimpanzees Carry Dead Infants: An Empirical Assessment of Existing Hypotheses*, 7 ROYAL SOC’Y OPEN SCI. 1 (2020) (observing that chimpanzees have death awareness, so behaviors such as carrying dead infants around lack an explanation

but it could be grief); Alexandra G. Rosati & Brian Hare, *Chimpanzees and Bonobos Exhibit Emotional Responses to Decision Outcomes*, PLOS ONE (May 29, 2013), <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0063058#s4> ("apes selectively attempted to switch their choices following undesired outcomes"); Alicia P. Melis et al., *Chimpanzees Coordinate in a Negotiation Game*, 30 EVOLUTION & HUM. BEHAV. 381 (2009) (describing experiments that show chimpanzees negotiate to resolve conflicts over resources); Alicia P. Melis & Michael Tomasello, *Chimpanzees' (Pan Troglodytes) Strategic Helping in a Collaborative Task*, 9 ROYAL SOC'Y BIOLOGY LETTERS 1 (2013) (describing an experiment where chimpanzees successfully completed tasks that required them to work together to get food); Brian Hare et al., *Chimpanzees Deceive a Human Competitor by Hiding*, 101 COGNITION 495 (2006) (describing an experiment where chimpanzees employed deception to outcompete humans to reach food).

29. Petitioners' Memorandum, *supra* note 7, at 4–5.

30. See, e.g., HERBERT S. TERRACE, WHY CHIMPANZEES CAN'T LEARN LANGUAGE AND ONLY HUMANS CAN (2019) (arguing that words are the cornerstone of language and that the 1973 Nim experiment demonstrated that Nim the chimpanzee could not learn words and thus could not acquire language; sign language was simply used to obtain rewards and thus lacked the cognition humans employ when using language).

31. Petitioners' Memorandum, *supra* note 7, at 16.

32. This is a personal, and very loose, translation from GUSTAVE FLAUBERT, MADAME BOVARY 201 (1901).

33. E.P EVANS, THE CRIMINAL PROSECUTION AND CAPITAL PUNISHMENT OF ANIMALS 18–19 (1906), <https://archive.org/details/criminalprosecut00evaniala/mode/2up>.

34. *Id.* at v–viii.

35. *Id.* at 150.

36. *Id.* at 186.

37. Felix Cohen, *Transcendental Nonsense and the Functional Approach*, 35 COLUM. L. REV. 809, 813 (1935).

38. John Dewey, *The Historic Background of Corporate Legal Personality*, 35 YALE L.J. 655, 656 (1926).

39. People ex rel. Nonhuman Rts. Project v. Lavery, 124 A.D.3d 149, 151 (N.Y. App. Div. 2014).

40. Nonhuman Rts. Project v. Lavery, 31 N.Y.3d 1054, 1057–58 (N.Y. 2018) (Fahey, J., concurring) (citations omitted).

41. Palila v. Hawaii Dep't of Land & Nat. Res., 852 F.2d 1106, 1107 (9th Cir. 1988).

42. Cetacean Cmty. v. Bush, 386 F.3d 1169 (9th Cir. 2004).

43. Naruto v. Slater, 888 F.3d 418, 420 (9th Cir. 2018) (emphasis added).

44. Christopher D. Stone, *Should Trees Have Standing? Toward Legal Rights for Natural Objects*, 45 S. CALIF. L. REV. 450 (1972); see also Christopher D. Stone, *Should Trees Have Standing? Revisited: How Far Will Law and Morals Reach? A Pluralist Perspective*, 59 S. CAL. L. REV. 1 (1985).

45. *Naruto*, 888 F.3d at 432 (Smith, J., concurring in part) (emphasis added).
46. *Dred Scott v. Sandford*, 60 U.S. 393, 426 (1857), superseded (1868).
47. *Id.* at 410.
48. *Id.*
49. *Kyllo v. United States*, 533 U.S. 27 (2001).
50. *Grady v. North Carolina*, 575 U.S. 306 (2015).
51. *Bostock v. Clayton Cnty*, 140 S. Ct. 1731, 1743 (2020).
52. *Id.* at 1755–56 (Alito, J., dissenting) (“[t]he Court’s opinion is like a pirate ship. It sails under a textualist flag, but what it actually represents is a theory of statutory interpretation that Justice Scalia excoriated—the theory that courts should ‘update’ old statutes so that they better reflect the current values of society”).
53. *Cetacean Cmty. v. Bush*, 386 F.3d 1169, 1176 (9th Cir. 2004).
54. See James Boyle, *The Anatomy of a Torts Class*, 34 AM. L. REV. 1003, 1056–57 (1985):

Courts are the Competent Institutions. This issue is uniquely suitable for the courts to deal with. Courts are the bodies that society has set up to deal with complex factual issues, to be responsive to changing circumstances, and yet to be objective. An issue like this, which needs all of these qualities for its resolution, must be left to the courts. In addition, this issue needs to be resolved by an institution, which can take outside expert advice and has a firm understanding of the changing moral consensus of our society. Courts are the only bodies which combine all of these abilities.

Courts are not the Competent Institution. This issue cannot be decided by the courts. Dealing as it does with issues of change. It must be left to the legislature, a body which reflects changing public opinion, which can bring in outside experts, and which is used to dealing with complex factual issues such as this. Regulation is the job of the legislature and the administrative branch; the courts should apply and not make the law; if they do otherwise they threaten the separation of powers. This issue, thus, cannot be seen outside the context of our entire institutional structure—a structure which should not be threatened for the sake of some “quick fix” judicial solution.
55. *Bowers v. Hardwick*, 478 U.S. 186 (1986), overruled by *Lawrence v. Texas*, 539 U.S. 558 (2003).
56. *Obergefell v. Hodges*, 576 U.S. 644 (2015).
57. *Id.* at 726 (Roberts, J., dissenting).
58. See David E. Bernstein, *Brandeis Brief Myths*, 15 GREEN BAG 2D 9 (2011) (discussing the importance of, and mythmaking behind, Brandeis’s famed brief).

CHAPTER 5

1. D. Scott Bennett, *Chimera and the Continuum of Humanity: Erasing the Line of Constitutional Personhood*, 55 EMORY LJ. 347, 348–49 (2006).
2. Henry T. Greely et al., *Thinking about the Human Neuron Mouse*, 7 AM. J. BIOETHICS 27 (2007). The report is a model of ethical clarity in thinking about these issues, not least because it focuses on practical concerns: the treatment of the animals, the

scientific and ethical benefits of the research and whether those can be achieved in other ways, the various reasons that people might object to the use of human tissues or cells, the types of reactions such research might engender in regulators, and so on. Bioethicists do a lot of that kind of work, carefully working out in the specific details of some research protocol and procedure the various ethical interests involved and administrative safeguards, or flat prohibitions, needed. They deserve our gratitude for it. Here, I am going to be focused on the more speculative, philosophical side of ethical work, but I do not want that to give a misleading impression of the field as a whole.

3. Sarah Taddeo, *Intraspecies Chimeras Produced in Laboratory Settings (1960–1975)*, EMBRYO PROJECT ENCYC. (Nov. 25, 2014), <https://embryo.asu.edu/pages/intraspecies-chimeras-produced-laboratory-settings-1960-1975> (emphasis in original). With intraspecies chimeras, an entity receives cells of a different genotype from another member of the *same* species. This can happen in bone marrow transplants or when one fetus also has cells derived from a sibling in the womb. I will be talking here only about *interspecies* chimeras, and specifically about human-nonhuman animal chimeras.

4. See Nidhi Subbaraman, *First Monkey–Human Embryos Reignite Debate over Hybrid Animals*, NATURE (Apr. 15, 2021), <https://www.nature.com/articles/d41586-021-01001-2>:

Scientists have successfully grown monkey embryos containing human cells for the first time—the latest milestone in a rapidly advancing field that has drawn ethical questions. In the work, published on 15 April in *Cell*, the team injected monkey embryos with human stem cells and watched them develop. They observed human and monkey cells divide and grow together in a dish, with at least 3 embryos surviving to 19 days after fertilization. . . . Researchers hope that some human-animal . . . chimaeras—could provide better models in which to test drugs, and be used to grow human organs for transplants. Members of this research team were the first to show in 2019 that they could grow monkey embryos in a dish for up to 20 days after fertilization. In 2017, they reported a series of other hybrids: pig embryos grown with human cells, cow embryos grown with human cells, and rat embryos grown with mouse cells.

5. Jamie Shreve, *The Other Stem Cell Debate*, N.Y. TIMES MAG. (Apr. 10, 2005), <https://www.nytimes.com/2005/04/10/magazine/the-other-stemcell-debate.html>.

6. K. Shankar & H.M. Mehendale, *Transgenic Animals*, in ENCYC. TOXICOLOGY (3d ed. 2014).

7. Monika Piotrowska, *What Does It Mean to Be 75% Pumpkin? The Units of Comparative Genomics*, 76 PHIL. SCI. 838, 838 (2009).

8. Shreve, *supra* note 5.

9. *Id.*

10. Kristin Hugo, *Exclusive: Whatever Happened to the Mouse with the Ear on Its Back?*, NEWSWEEK (Sep. 16, 2017), <https://www.newsweek.com/tissue-surgeon-ear-mouse-human-organs-transplant-cell-phones-666082>.

11. Greely et al., *supra* note 2, at 36.

12. *Id.* at 27.

13. S. 1373, 109th Cong. § 2 (2005). Note, the bill was first introduced as S. 659, 109th Cong. § 2 (2005).
14. George W. Bush, *State of the Union Address by the President*, WHITE HOUSE (Jan. 31, 2006), <https://georgewbush-whitehouse.archives.gov/stateoftheunion/2006> (please note that this web page is “frozen” for archival purposes, so none of the links that appear on the web page will work).
15. U.S. Patent No. 4,736,866.
16. Greely et al., *supra* note 2, at 32.
17. It may be disciplinary prejudice, but I found the articles by Greely et al. to be both more practical about regulatory choices and more wide-ranging in their discussion of the ethical issues involved.
18. See, e.g., César Palacios-González, *Human Dignity and the Creation of Human-Nonhuman Chimeras*, 18 MED., HEALTH CARE, & PHIL. 487 (2015); Phillip Karpowicz et al., *Developing Human-Nonhuman Chimeras in Human Stem Cell Research: Ethical Issues and Boundaries*, 15 KENNEDY INST. ETHICS J. 107 (2005); Phillip Karpowicz et al., *It Is Ethical to Transplant Human Stem Cells into Nonhuman Embryos*, 10 NATURE MED. 331 (2004); Josephine Johnston & Christopher Eliot, “*Chimeras*” and Human Dignity, 3 AM. J. BIOETHICS W6 (2003); CHIMERA’S CHILDREN: ETHICAL, PHILOSOPHICAL AND RELIGIOUS PERSPECTIVES ON HUMAN-NONHUMAN EXPERIMENTATION (David Albert Jones & Calum MacKellar eds., 2012).
19. For an accessible overview of the philosophical debates about moral status, see Agnieszka Jaworska & Julie Tannenbaum, *The Grounds for Moral Status*, STAN. ENCYCL. PHIL. (Edward N. Zalta & Uri Nodelman eds., Mar. 3, 2021), <https://plato.stanford.edu/archives/spr2023/entries/grounds-moral-status>. For specific discussions of different conceptions of what confers moral status, see, e.g., IMMANUEL KANT, GROUNDWORK FOR THE METAPHYSICS OF MORALS (Mary Gregor & Jens Timmermann eds. trans., Cambridge Univ. Press 2012) (1785); Warren Quinn, *Abortion: Identity and Loss*, 13 PHIL. & PUB. AFFS. 24 (1984); JEFF MACHMAHON, THE ETHICS OF KILLING: PROBLEMS AT THE MARGINS OF LIFE (2012); Michael Tooley, *Abortion and Infanticide*, 2 PHIL. & PUB. AFFS. 37 (1972); Sarah Buss, *The Value of Humanity*, 109 J. PHIL. 341 (2012); L. NANDI THEUNISSEN, THE VALUE OF HUMANITY (2020); Joel Feinberg, *Abortion*, in MATTERS OF LIFE AND DEATH 256 (Tom L. Beauchamp & Tom Regan eds., 1980); Agnieszka Jaworska, *Caring and Full Moral Standing*, 117 ETHICS 460 (2007).
20. Robert Streiffer, *In Defense of the Moral Relevance of Species Boundaries*, 3 AM. J. BIOETHICS 37, 38 (2003).
21. Leon Kass, *The Wisdom of Repugnance*, NEW REPUBLIC, June 2, 1997, at 17.
22. *Id.*
23. See FYODOR DOSTOEVSKY, CRIME AND PUNISHMENT (1866).
24. Andrew Siegel, *The Moral Insignificance of Crossing Species Boundaries*, 3 AM. J. BIOETHICS 33, 33 (2003).
25. Palacios-González, *supra* note 18, at 490.

26. See Jeff McMahan, *Cognitive Disability, Misfortune, and Justice*, 25 PHIL. & PUB. AFFS. 3, 30 (1996). Perhaps to his credit, McMahan is willing to say so clearly and specifically rather than to obfuscate.

To my mind, the most plausible general view is that there are certain properties and capacities that give their possessor an inherent worth that demands respect. It is the possession of these properties and capacities that makes an individual one's moral equal and thus brings him or her within the sphere of justice. There are different accounts of what the relevant properties and capacities are, though there is general agreement that they are psychological rather than physical in nature . . . if the properties are intrinsic rather than relational, and in particular if they are psychological rather than physical, and if animals (or at least animals other than the great apes) are excluded from the sphere of justice because they do not possess these properties, *then it seems to follow that human beings with comparable psychological properties and capacities must be excluded as well* [emphasis added].

I recommend the response of Eva Feder Kittay, a philosopher and ethicist who is herself the mother of a daughter with mental disabilities. Eva Feder Kittay, *At the Margins of Moral Personhood*, 116 ETHICS 100 (2005). The emotion and compassion with which she sees the issue brilliantly animates but does not overwhelm the careful ethical argument, rather than simply prompting an outpouring of outrage.

27. Karpowicz et al., *supra* note 18, at 118.

28. Loving v. Virginia, 388 U.S. 1, 3 (1967) (quoting the Caroline County Circuit Court but providing no formal citation to a reporter; Caroline County has only published decisions from 1995 and unpublished decisions from 2002 available for search on its website at <https://www.vacourts.gov/online/home.html>).

29. Loving v. Virginia, 388 U.S. 1, 7 (1967).

30. Palacios-González, *supra* note 18, at 490.

31. This possibility has been the basis for much fascinating science fiction. See, e.g., DAVID BRIN, THE UPLIFT WAR (1987). I am not alone in thinking that the moral issues around Brin's uplift novels are given new salience by AI.

COWEN: What in science fiction do you feel has risen the most in status for you?

HOFFMAN: Oh, for me.

COWEN: Not in the world. We don't know yet.

HOFFMAN: Yes. We don't know yet.

COWEN: You think, "Oh, this was really important." Vernor Vinge or . . .

HOFFMAN: Well, this is going to seem maybe like a strange answer to you, but I've been rereading David Brin's Uplift series very carefully because the theory of, "How should we create other kinds of intelligences, and what should that theory be, and what should be our shepherding and governance function and symbiosis?" is a question that we have to think about over time. He went straight at this in a biological sense, but it's the same thing, just a different substrate with the Uplift series. I've recently reread the entire Uplift series.

See Conversations with Tyler, *Reid Hoffman on the Possibilities of AI* (June 28, 2023), <https://conversationswithtyler.com/episodes/reid-hoffman-2>.

32. Streiffer, *supra* note 20, at 38.

33. It should be noted that variants of this argument can be used to attack as well as to apologize for the existing social order. The Levellers' subversive ditty "When Adam delv'd and Eve span / Who then was the gentleman?" is both a good example of such an attack and my favorite piece of rhyming natural law. If there was no social hierarchy in the Garden of Eden, the divine state of nature, it follows that social stratification is something human created, *after* the Fall. It thus must be a result of the entry of evil into the world. The nobility are the work of the devil!

34. S. 1373, 109th Cong. § 2 (2005). Note, the bill was first introduced as S. 659, 109th Cong. § 2 (2005).

35. EDMUND BURKE, REFLECTIONS ON THE REVOLUTION IN FRANCE 55 (1790).

36. EDMUND BURKE, *A Letter to a Member of the National Assembly* (1791), in THE WORKS OF THE RIGHT HONOURABLE EDMUND BURKE, VOL. IV (Project Gutenberg EBook, 2005), https://www.gutenberg.org/files/15700/15700-h/15700-h.htm#MEMBER_OF_THE_NATIONAL_ASSEMBLY.

37. Jason Scott Robert & Francoise Baylis, *Crossing Species Boundaries*, 3 AM. J. BIOETHICS 1, 9 (2003).

38. *Id.*

39. Streiffer, *supra* note 20, at 37.

40. *Id.*

41. David Bowie, *Rebel Rebel*, on DIAMOND DOGS (RCA Records, 1974).

42. Leon Kass, *The Wisdom of Repugnance*, NEW REPUBLIC, June 2, 1997 (emphasis added).

43. Karpowicz et al., *supra* note 18, at 120 (emphasis added).

44. See Robert Streiffer, *Human/Non-Human Chimeras*, STAN. ENCYC. PHIL. (Edward N. Zalta ed., Apr. 5, 2019), <https://plato.stanford.edu/entries/chimeras> (discussing the conflict between the "Millian"/utilitarian view and the Kantian view in the "The Moral Status Framework" section); see also Robert Streiffer, *At the Edge of Humanity: Human Stem Cells, Chimeras, and Moral Status*, 15 KENNEDY INST. ETHICS J. 347 (containing Streiffer's formal arguments for his moral status framework).

45. I owe the phrase to the excellent science fiction of Joel Shepherd.

46. FRANS DE WAAL, ARE WE SMART ENOUGH TO KNOW HOW SMART ANIMALS ARE? 5 (2016).

47. Karpowicz et al., *Developing Human-Nonhuman Chimeras* *supra* note 18, at 121–22.

48. BERNARD WILLIAMS, *Persons, Character and Morality*, in MORAL LUCK: PHILOSOPHICAL PAPERS 1973–1980, at 1, 18 (1981).

49. Sophie-Grace Chappel & Nicholas Smyth, *Bernard Williams*, STAN. ENCYC. PHIL. (Edward N. Zalta & Uri Nodelman eds., Jan. 28, 2023), <https://plato.stanford.edu/entries/williams-bernard>.

50. Buck v. Bell, 274 U.S. 200, 207 (1927).

51. JAMES BOYLE, SHAMANS, SOFTWARE AND SPLEENS (1996).
52. James Boyle, *Endowed by Their Creator? The Future of Constitutional Personhood*, in CONSTITUTION 3.0: FREEDOM AND TECHNOLOGICAL CHANGE (Jeff Rosen & Benjamin Wittes eds., 2013) (the book was published in 2013, but the website version was published in 2011).
53. Karpowicz et al., *supra* note 18, at 121–22.
54. Those who believe that wrongly projecting humanity onto potentially manipulative Artificial Intelligences carries with it the substantial risk of killing off the entire human species will, understandably, disagree. Fair point.

CONCLUSION

1. PETER SINGER, ANIMAL LIBERATION (1975); Matthew D. Adler, *Future Generations: A Prioritarian View*, 77 GEO. WASH. L. REV. 1478 (2009).
2. ADAM SMITH, THE THEORY OF THE MORAL SENTIMENTS (Knud Haakonssen ed., Cambridge Univ. Press 2002) (1759); DAVID HUME, A TREATISE OF HUMAN NATURE (David Fate Norton & Mary J. Norton eds., Oxford Philosophical Texts, 2006) (1739–40).
3. JEREMY BENTHAM, AN INTRODUCTION TO THE PRINCIPLES OF MORALS AND LEGISLATION, ch. xvii, at 351 (1789), <http://www.koeblergerhard.de/Fontes/BenthamJeremyMoralsandLegislation1789.pdf>.
4. Stephen Hawking et al., *Stephen Hawking: Transcendence Looks at the Implications of Artificial Intelligence: But Are We Taking AI Seriously Enough?*, INDEPENDENT (May 1, 2014), <https://www.independent.co.uk/news/science/stephen-hawking-transcendence-looks-at-the-implications-of-artificial-intelligence-but-are-we-taking-ai-seriously-enough-9313474.html> (emphasis added).
5. 2022 Expert Survey on Progress in AI, AI IMPACTS WIKI (Aug. 4, 2022), https://wiki.aiimpacts.org/doku.php?id=ai_timelines:predictions_of_human-level_ai_timelines:ai_timeline_surveys:2022_expert_survey_on_progress_in_ai.
6. See also, James Somers, *The Pastry A.I. That Learned to Fight Cancer*, NEW YORKER (Mar. 18, 2021), <https://www.newyorker.com/tech/annals-of-technology/the-pastry-ai-that-learned-to-fight-cancer> (discussing how an AI developed for distinguishing among types of pastries is also being used to help identify cancer cells); Stephen Ornes, *The Unpredictable Abilities Emerging from Large AI Models*, QUANTA MAG. (Mar. 16, 2023), <https://www.quantamagazine.org/the-unpredictable-abilities-emerging-from-large-ai-models-20230316> (discussing how ChatGPT can do multiplication, write usable computer code, and decode movies based on emojis).
7. Kevin Roose, *Why an Octopus-Like Creature Has Come to Symbolize the State of A.I.*, N.Y. TIMES (May 30, 2023), <https://www.nytimes.com/2023/05/30/technology/shoggoth-meme-ai.html>.
8. B.F. SKINNER, CONTINGENCIES OF REINFORCEMENT 260 (Copley Publ'g Grp. 2013) (1969).
9. STUART RUSSELL & PETER NORVIG, ARTIFICIAL INTELLIGENCE: A MODERN APPROACH 3 (3d ed. 2010).

10. Stephen Wolfram, *What Is ChatGPT Doing and Why Does It Work?*, STEPHEN WOLFRAM: WRITINGS (Feb. 14, 2023), <https://writings.stephenwolfram.com/2023/02/what-is-chatgpt-doing-and-why-does-it-work>.
11. Alexander Pope, *An Essay on Man: Epistle II*, POETRY FOUND, <https://www.poetryfoundation.org/poems/44900/an-essay-on-man-epistle-ii> (last visited July 19, 2023).
12. Trustees of Dartmouth College v. Woodward, 17 U.S. 518, 636 (1819).
13. First Nat'l Bank of Boston v. Bellotti, 435 U.S. 675, 822 (1978) (Rehnquist, J., dissenting).
14. *Id.* at 822–25 (Rehnquist, J., dissenting).
15. HOWARD J. GRAHAM, “*Buidled Better Than They Knew*: The Framers, the Railroads and the Fourteenth Amendment, in EVERYMAN'S CONSTITUTION 443 (Wis. Hist. Soc'y 2013) (1968).
16. Citizens United v. FEC, 558 U.S. 310, 466 (2010) (Stevens, J., dissenting).
17. SMITH, *supra* note 2, at 11.
18. Jamal Greene, *Originalism's Race Problem*, 88 DENVER L. REV. 517 (2011).
19. See, e.g., James Boyle, *Endowed by Their Creator? The Future of Constitutional Personhood*, in CONSTITUTION 3.0: FREEDOM AND TECHNOLOGICAL CHANGE (Jeff Rosen & Benjamin Wittes eds., 2013) (the discussion of “Vanna”).
20. See, e.g., Mita Giacomini, *A Change of Heart and a Change of Mind? Technology and the Redefinition of Death in 1968*, 44 SOC. SCI. MED. 1465 (1997) (evaluating the role of life support, EEG, and organ transplantation in influencing the change in the definition of death); Amir Halevy & Baruch Brody, *Brain Death: Reconciling Definitions, Criteria, and Tests*, 119 ANNALS INTERNAL MED. 449 (discussing the limitations of diagnostic practices in determining whether brain activity has truly ceased); Steven Goldberg, *The Changing Face of Death: Computers, Consciousness, and Nancy Cruzan*, 43 STAN. L. REV. 659 (1991) (noting how developments in technology precipitated changes in the legal definition of death and how new developments in artificial intelligence may lead to new changes in the definition of death); John P. Lizza, *Defining Death: Beyond Biology*, 55 DIAMETROS 1 (2018) (defending brain death as death by arguing that artificially sustained bodies that are brain-dead are not human persons); but see D. Alan Shewmon, *Brain Death: Can It Be Resuscitated*, 39 HASTINGS CTR. REP. 18, 19 (2009) (noting that the general public and many in the medical profession consider brain-dead patients to be “as good as dead” or “better off dead” rather than “really dead”).
21. See, e.g., Robert D. Truog & James C. Fackler, *Rethinking Brain Death*, 20 CRITICAL CARE MED. 1705 (1992) (arguing for a revision of the definition of brain death to a “higher brain” standard instead of a “whole brain” standard); but see Robert D. Truog, *Is It Time to Abandon Brain Death?*, 27 HASTINGS CTR. REP. 29 (1997) (arguing in favor of abandoning the concept of brain death altogether and uncoupling it from organ transplantation); Robert M. Taylor, *Reexamining the Definition and Criteria of Death*, 17 SEMINARS NEUROLOGY 265, 265 (1997) (arguing that “[t]he best definition of death is ‘the event that separates the process of dying from the process of

disintegration' and the proper criterion of death in human beings is 'the permanent cessation of the circulation of blood'" but that changing to such a definition would be politically problematic).

22. Lydia Wheeler, *Fetal Rights Laws' Impact Extends from Abortion to HOV Lanes*, BLOOMBERG (July 27, 2022), <https://news.bloomberglaw.com/us-law-week/fetal-rights-laws-impact-extends-from-abortion-to-hov-lanes>.

23. U.S. CONST. amend. XIV § 1.

24. Advocates, including the natural law scholar John Finnis, have used an argument similar to the ones Tommy's lawyers employed: that by various criminal, tort, and other remedies, the substantive law of the states *already* recognized fetuses as legal persons when the Fourteenth Amendment was enacted. Brief for Scholars of Jurisprudence John M. Finnis & Robert P. George as Amici Curiae Supporting Petitioners, *Dobbs v. Jackson Women's Health Org.*, 142 S. Ct. 2228 (2022) (No. 19-1392); see also Gregory J. Roden, *Unborn Children as Constitutional Persons*, 25 ISSUES L. MED. 185 (2010). To me, the historical record clearly indicates the opposite. Laws refer to or regulate many subjects without transmuting those subjects into legal persons and congressional practice indicates a specific decision to exclude the unborn in multiple contexts, some of them directly relevant to the Fourteenth Amendment. See, e.g., Michael L. Rosin, *Congress Has Never Considered Fetuses Persons within the Meaning of the 14th Amendment*, SLATE (June 9, 2022), <https://slate.com/news-and-politics/2022/06/gop-abortion-constitution-fetuses-legal-persons.html>.

25. Howard J. Graham, *The "Conspiracy Theory" of the Fourteenth Amendment*, 47 YALE L.J. 371, 378 (1938) (quoting Oral Argument of Roscoe Conkling, in SAN MATEO CASE, ARGUMENTS AND DECISIONS 33 [Stan. Univ. Library ed.]); see also Ralph Waldo Emerson, *The Problem*, POETS.ORG, <https://poets.org/poem/problem> (last visited July 20, 2023) ("[h]e bulidde better than he knew").

26. See Natasha N. Aljalian, Note, *Fourteenth Amendment Personhood: Fact or Fiction?*, 73 ST. JOHN'S L. REV. 495 (1999) (arguing that if corporations meet the constitutional definition of personhood then fetuses should as well); see also Michael Stokes Paulsen, *The Plausibility of Personhood*, 74 OHIO ST. L.J. 13, 23 n.34 (2013) (noting the argument from advocates that runs from "irony or hypocrisy" that "[i]f even corporations can be treated by the law as legal 'persons,' then surely unborn babies should be treated as legal persons, too").

27. See, e.g., LEE EPSTEIN & JACK KNIGHT, *THE CHOICES JUSTICES MAKE* (1997) (giving a strategic account of how Justices on the Supreme Court behave to achieve their policy preferences); Jack Knight & Lee Epstein, *The Norm of Stare Decisis*, 40 AM. J. POL. SCI. 1018 (1996) (describing how Justices account for the norm of stare decisis in their decision-making); Thomas J. Miceli & Metin M. Coşgel, *Reputation and Judicial Decision-Making*, 23 J. ECON. BEHAV. & ORG. 31 (1994) (describing the role of reputation in judicial decision-making); Gregory C. Sisk et al., *Charting the Influences on the Judicial Mind: An Empirical Study of Judicial Reasoning*, 73 N.Y.U. L. REV. 1377 (1998) (analyzing a mix of legal and extra-legal factors on judicial decision-making); Allison P. Harris & Maya Sen, *Bias and Judging*, 22 ANN. REV. POL. SCI. 241 (2019) (discussing the

role of ideology in judicial decision-making relative to bias and ideology's greater importance).

28. See Michael J. Gerhardt, *Super Precedent*, 90 MINN. L. REV. 1204, 1205 (2006) (arguing for the existence of superprecedents in agreement with Chief Justice Roberts); see also Dobbs v. Jackson Women's Health Org., 142 S. Ct. 2228, 2262–78 (2022) (analyzing whether stare decisis protects *Roe v. Wade* and *Planned Parenthood v. Casey* from being overruled); *Confirmation Hearing on the Nomination of Hon. Brett M. Kavanaugh to Be an Associate Justice of the Supreme Court of the United States: Before the S. Comm. on the Judiciary*, 115th Cong. 127 (2018) (statement of Brett Kavanaugh, Judge, United States Court of Appeals for the District of Columbia) (testifying that *Roe v. Wade* was settled precedent).

29. See Biden v. Nebraska, 143 S. Ct. 2355, 2373–75 (2023) (deciding that student loan cancellation was a major question that Congress would have reserved for itself).

30. Mark A. Lemley, *The Imperial Supreme Court*, 136 HARV. L. REV. F. 97, 97 (2022) (“the Court has begun to implement the policy preferences of its conservative majority in a new and troubling way: by simultaneously stripping power from every political entity except the Supreme Court itself”).

31. Citizens United v. FEC, 558 U.S. 310, 466 (2010) (Stevens, J., dissenting).

32. See generally Paul F. Christiano, *Induction; Or, The Rules and Etiquette of Reference Class Tennis*, LESSWRONG (Mar. 3, 2018), <https://www.lesswrong.com/posts/PXRxH4C6nKMwocBit/induction-or-the-rules-and-etiquette-of-reference-class>.

33. Andrew Lakoff, *Vaccine Politics and the Management of Public Reason*, 27 PUB. CULTURE 419, 419–20 (2015).

34. *Id.*

35. Rund Abdelfatah, “*Throughline*” Traces Evangelicals’ History on the Abortion Issue, NPR MORNING EDITION (June 20, 2019), <https://text.npr.org/734303135>.

36. *Id.*

37. As someone who has presented competing “ideal types” of arguments before, I know that one of the dangers in doing so is that lazy commentators will select one of the extreme views and present it as if it were my own position. I can hardly prevent that process but serve advance notice of my objection!

INDEX

- Abdelfatah, Rund, 314n35
Abelson, Hal, 80
Abney, Keith, 287n2
Abortion, and personhood debate, 21, 24, 29, 37, 188, 200, 215, 257–269, 308n19, 313n22, 313n24, 313n26, 314n35
Acharya, Sourya, 296n100
Adams, Ansel, 295n94
Adler, Matt, 218, 237, 311n1
Alfinito, E., 292n54
Alito, Samuel, 306n52
Aljalian, Natasha N., 313n26
Allen-Hermanson, Sean, 289
AlphaGo, 81, 291n37
 AlphaGo Zero, 81, 290n16, 291n39
American Sign Language, 25, 167, 177, 244
Animals, nonhuman
 appropriate decision-making bodies for rights' claims of, 188–193
 Darwin's claims about mental states in, 168–170
 habeas corpus claims of, 163, 176, 181–182, 199, 225, 303n7
 history of legal personalization of, 178–180
 insights derived from animal lawsuits, 183–193
 judicial methodologies in nonhuman rights cases, 185–188
legal rightsholder argument for, 178–183
legal strategies on behalf of, 165–183
mental capacities of, vs. humans, 168–178
Nonhuman Rights Project litigation for, 163, 165, 166–183
People ex rel. Nonhuman Rts. Project v. Lavery, 303n7, 305n39
personhood claims on behalf of, 163–193
qualities-of-mind argument for, 166–178
rats of Autun, 178–179
Tommy, legal claims of, 163–193
who speaks for the voiceless?, 183–184
Anthropomorphism, 6, 29, 33–35, 111, 120, 171, 179
and empathy, 33–35
pathetic fallacy's connection to, 34
Skinner and Morgenbesser discussing, 91 (*see also* Skinner, B. F.)
Arbus, Diane, 42
Aristotle, 25, 111, 117, 166–168, 171, 219, 220, 230, 244, 303n11
Art, 18–19, 112–118
and debate about machine consciousness, 18–19, 23, 29, 35
definitions of, 26, 112–114
possibility of machine-generated, 112–114, 118

- Art (cont.)
 redefinition of human artistry, 112–118, 242, 245, 272, 275
- Artificial Intelligence
 art created by, 112–114, 118
 copyright status of works created by, 112, 284–285n29, 295–296n94
 dangers of, 35, 66, 99–108, 239–241, 275, 311n54
 defined, 3, 289n8
 and definition of humanity, 3–4, 12, 25, 57, 110–118, 167, 242–243, 272
 EU draft proposal about, 132
 fictional explorations of, 31, 41–57
 Hal thought experiment about, 9–14
 and human brain capabilities, 78–80
 mistaken perceptions of consciousness in large language models, 1–5
 optimistic predictions about, 63–66, 82–83
 Platt's Law of AI prediction, 77, 79
 possibility of consciousness in, 83–97
 possible conservative and liberal perspectives on, 267–270
 possible personhood of, 1–14, 18–32, 83–97
 Searle's Chinese Room and AI “consciousness,” 84–94
 and Singularity (*see* Singularity, the)
 speed of development, estimating, 66–83, 101
 technical possibility of, 17–18, 66–83, 101
 terminology used to discuss, 3, 289n8
- Barnett, Mathew, 290n25
- Barthes, Roland, 296n95
- Bates, Laura, 288n26
- Bates, Mary, 303n17
- Baumol, William, 114, 296n96
- Baylis, Francoise, 214, 215, 310n37
- Beard, Charles, 150, 248, 300n34
- Beard, Mary, 248, 300n29
- Beatty, Jack, 301n38
- Beauchamp, Tom L., 308n19
- Bender, Emily, 3
- Ben-Menahem, Yemima, 294n65
- Bennett, D. Scott, 306n1
- Bentham, Jeremy, 166, 218, 238, 302n4, 303n10, 311n3
- Berman, Mitchell, 301n49
- Berners-Lee, Tim, 68
- Bernstein, David E., 306n58
- Biden v. Nebraska*, 315n29
- Bill of Rights, 23
 First Amendment, 22, 134, 145, 146, 300n28
 Fourth Amendment, 186
 Second Amendment, 186, 269
- Birhane, Abeba, 282n13
- Black, Hugo, 158
- Blade Runner*, 9, 19, 30, 35, 41–50, 53–56, 239, 242, 262, 287n7
- Batty, Roy, 48, 52, 54, 56–57
- Deckard, Rick, 42–47, 52, 56–57
- Bloom, Paul, 38, 287n12
- Bostock v. Clayton Cnty.*, 306n51–52
- Bostrom, Nick, 66, 78–79, 82, 102–105, 107–108, 116, 275, 289n11, 294n74, 294nn81–83
- Bowers v. Hardwick*, 306n55
- Bowie, David, 310n41
- Boyle, James, 281n6, 283n17, 285n29, 301n49, 306n54, 311nn51–52, 312n19
- Bradbury, Ray, 50
- Brandeis, Louis, 192, 306n58
- Brent, Michael, 283n18
- Brin, David, 309n31
- Britten, Roy J., 283n25
- Brod, Max, 45
- Brody, Baruch, 312n20
- Brooks, Rodney, 75–77, 291n27, 291nn29–30
- Brown, Arthur E., 304n28
- Brown, Teneille R., 285n29, 298n8

- Brownback, Sam, 202–203, 212–213, 215, 229
- Bryson, Joanna J., 282nn14–15
- Buccafusco, Chris, 277
- Buckner, Randy L., 288n27
- Buck v. Bell*, 310n50
- Burckhardt, John Lewis, 288n31
- Burger, Warren E., 146, 300n28
- Burke, Edmund, 212, 310nn35–36
- Burkean philosophy, 66, 212, 214, 255
- Bush, George W., 200, 203, 308n14
- Bush v. Gore*, 156
- Buss, Sarah, 308n19
- Butler, Samuel, 56–63, 86–88, 92, 172–173, 218, 288n30 (chap. 1), 288n2 (chap. 2), 289n3, 289n6
- consciousness in vegetables and objects, 61–62
- critiques of Victorian society, 59–60
- Erewhon*, 9, 59, 63, 92, 288n1, 304n22
- human vs. machine consciousness, 60
- possibility of machine “evolution,” 59–63
- Butlerian Jihad, 60, 99, 241
- Butlin, Patrick, 293n64, 294n68
- Calo, Ryan, 285n29
- Čapek, Karel, 18, 99, 286n30
- Caplan, Bryan, 75, 291n26
- Case of Sutton's Hospital*, 135, 299n9
- Cetacean Cnty. v. Bush*, 183–184, 187, 305n42, 306n53
- Chappel, Sophie-Grace, 310n49
- Charles, Guy-Uriel, 277
- ChatGPT, 3–5, 29–30, 63, 71, 81, 88, 102, 110–121, 167, 240, 242–244, 249, 270, 274, 281n8, 283n20, 311n6, 312n10
- Chimeras, 5, 8, 14, 16, 31, 164, 170, 183–185, 195–198, 202–205, 207–210, 212–215, 220, 226, 253, 255, 258–259, 261–262, 267, 270, 307n3, 308n18, 310n44, 310n47. *See also*
- Genetically engineered entities; Hybrids; Transgenic species
- Chimpy thought experiment, 14–17, 20–21, 56, 196–202, 231–236, 257, 263, 270–271, 284n28
- Chinese Room, 5, 84, 86–87, 89, 90, 94, 98, 110–111, 115, 120, 132, 242, 262, 272, 274
- Chomsky, Noam, 296n101
- Christiano, Paul F., 314n32
- Churchill, Winston, 286n29
- Citizens United v. FEC*, 133, 144–145, 157–159, 188, 247, 285n29, 298nn7–8, 299n22, 300n23, 300n28, 301n48, 301n52, 312n16, 314n31
- Civil War Amendments, 147
- Fifteenth Amendment, 7, 147
- Fourteenth Amendment, 7, 13, 143–159, 229, 248, 257–260, 300n29, 300n31, 300n35, 301n35, 301n41, 301n44, 312n15, 313nn23–26
- Thirteenth Amendment, 7, 13, 16, 147
- Clark, Betsy, 39, 40, 41, 287n13
- Clarke, James Freeman, 287n15
- Clifford, Catherine, 294n75
- Climate change, 14, 170, 268
- Clinton, Hillary, 144
- Coffee, John C., Jr., 300n24
- Cohen, Felix, 26, 132–133, 137–138, 140, 142, 165, 180–181, 251, 286n36, 298n1, 299n13, 303n8, 305n37
- Coke, Edward, 135–136
- Colangelo, P. Zak, 285n29
- Collins, Charles Wallace, 153, 301n41
- Collins, Eli, 281n3
- Conkling, Roscoe, 147–151, 153, 155–156, 158, 228, 260, 300n29, 313n25
- Consciousness
- computational functional theories of, 96–97, 293n65
- contemporary neuroscientific theories of, 95–98

- Consciousness (cont.)
 global neuronal workspace theory of, 96–97, 120
 as an illusion, 26, 56, 94–95
 impossibility of species-independent account of, 216–221, 234, 256
 integrated information theory of, 95–96, 120
 large language models' lack of, 26
 panpsychist theories of, 62, 289n5
 possibility of, in machines, 83–97
 recurrent processing theory, 97
 skepticism about existence in humans, 26, 56, 94–95
- Constantin, Sarah, 290n16
- Corporate personality
 associational theory of, 139
 and *Citizens United* (see *Citizens United v. FEC*)
 competing theories of, 136–141
 and Conkling, Roscoe (see Conkling, Roscoe)
 and constitutional rights, 143–157
 equal protection arguments for and against, 157–159
 equal protection rights “without argument or discussion,” 145–157
 and First Amendment claims, 145–146
 and Fourteenth Amendment claims, 143–144, 147–157
 legal fiction theory of, 137–139
 nexus of contracts theory of, 140–141
 as preview of AI debate, 141–143, 159–161
 real entity theory of, 136–137
 Rehnquist's critique of constitutionalization, 145–146
 and Story, Joseph, 131–133 (see also Story, Joseph)
- Coşgel, Metin M., 313n27
- Cotra, Ajeya, 79, 291n36
- Cowen, Tyler, 72–73, 109, 290n19, 290n21, 291n42, 295n92, 297n103, 298n111 (chap. 2), 298n3 (chap. 3), 309n31
- Cox, John Woodrow, 281n2
- Crichton, Michael, 53, 288n28
- Cuthbertson, Anthony, 290n18
- Darling, Kate, 34, 164, 285n29, 287n2, 302n3, 303n9
- Darwin, Charles, 168–169, 171, 173, 175, 178, 288n2, 303nn15–16
- Darwin, Erasmus, 288n2
- Davis, J. Bancroft, 152–153, 155
- Dawson, Michelle, 109, 291n42, 295n92
- Death, changing definitions of, 257–259, 312nn20–21
- Declaration of Independence, 17, 20, 185–186, 269
- Definitions of phenomena
 critiques by Hobbes, Wittgenstein, and Cohen, 26
 philosophical dangers of, 26–28
- DeGrazia, David, 302n4
- Delvaux, Mady, 132–133
- Dennett, Daniel, 87–88, 91, 218, 284n28, 292n48, 302n4
- Descartes, René, 86
- de Waal, Frans, 61, 170–173, 178, 219–220, 249, 303n17, 304nn18–21, 310n46
- Dewey, John, 137–138, 140, 142, 180–181, 251, 299n14, 305n38
- Dicey, A. V., 136, 299n11
- Dick, Philip K., 35, 39, 41, 45, 47, 49, 121, 239, 287n8, 287nn20–21
- Do Androids Dream of Electric Sheep?* (Dick), 41–42, 46, 49, 55–56, 108, 160, 239
- explorations of empathy in, 43–46
- Voight-Kampff Test in (see Voight-Kampff Test)
- Dobbs v. Jackson Women's Health Organization*, 155, 188, 313n24, 314n28

- Doctorow, Cory, 42, 287n19
Dostoevsky, Fyodor, 308n23
Douglas, William O., 158, 301n47
Dred Scott v. Sandford, 14, 185,
 306nn46–48
DuBois, W. E. B., 273
- Edison, Thomas, 72
Einstein, Albert, 68
Eliot, Christopher, 308n18
Emerson, Ralph Waldo, 149, 229,
 300n33, 313n25
Empathy, 30, 33, 39, 84, 193, 246
 criticisms of, 37–38
 as one basis for morality, 35–41
 use by anti-slavery movement, 39–41
 and Voight-Kampff Test, 43–44
Environmental precautionary principle,
 211, 213
Epstein, Lee, 313n21
Erewhon (Butler), 9, 59, 63, 92, 288n1,
 304n22
Etzioni, Oren, 73
Evans, E. P., 179, 305nn33–36
Evolution, 60, 106, 172, 206, 242, 249
 denunciations of, 60–61, 88
 process of, 9, 13, 60–61, 79, 87–88, 92,
 168, 173, 218, 220
 theory of, 25, 28, 168–169, 214, 232
Exponential growth of technology, 65,
 67–68, 70–73, 75, 78, 100–101
 and Moore’s law, 67, 70, 79, 100
- Facebook, 102
Fackler, James C., 312n21
Fahey, Eugene M., 182, 305n40
Feder Kittay, Eva, 309n26
Feinberg, Joel, 308n19
Feingold, Russ, 144–145
Ferrara, Ralph C., 300n24
Feynman, Richard, 68
Finkel, Elizabeth, 293nn61–62
Finnigan, Bronwyn, 303n13
- Finnis, John, 313n24
First National Bank of Boston v. Bellotti,
 145, 247, 300n26, 301n45, 312n13
Fischetti, Mark, 291n32
Fisher, Marc, 281n2
Flaubert, Gustave, 305n56
Forbes, 281n8
Fost, Norman, 200
Frankish, Keith, 293n56
Friedman, Milton, 269
Froomkin, Michael, 285n29
- Gallagher, Shaun, 296n101
Garreau, Joel, 287n1
Garrett, Brandon, 298n108, 298n110
General AI, 3–4, 30, 64–66, 71–77, 80,
 81, 82, 83, 99, 100–101, 105, 120,
 125–126, 164, 239–242, 275, 289n8
 arrival of, 65, 67, 72, 74–75, 77–79,
 82, 101, 239
 Literal and Rogue AI, 104–107, 127,
 270
 and reinforcement learning, 70–71,
 80–81
Genetically engineered entities
 and *Blade Runner*, 46, 49–57
 and blurring of species line, 203, 214–
 215, 254–255, 269
 and Chimpy thought experiment (*see*
 Chimpy thought experiment)
 and definitions of humanity, 198–202
 DNA percentage similarity to humans,
 198–199
 and human dignity, 198–202, 205
 legislative proposals to ban, 202–204
 portrayal concerns about human
 likeness, 200–201, 254
possible personhood claims on behalf
 of, 6, 14–17
potential-based definitions and,
 201–202
procreation-based definitions of
 humanity and, 200

- Genetically engineered entities (cont.)
 provenance of cells and, 199, 256
 and use in medical research, 195–215
- George, Robert P., 313n24
- Gerhardt, Michael J., 314n28
- Gersen, Jeannie Suk, 288n26
- Gervais, Daniel, 285n29
- Ghahramani, Zoubin, 281n3
- Giacomini, Mita, 312n20
- Goff, Philip, 289n5
- Goldberg, Steven, 312n20
- Good, Irving John, 99, 294n71
- Google, 1, 4, 9, 11–14, 142, 243, 261,
 281n1, 281n3, 281n10, 290n17
- Bard, 3, 281n10
- DeepMind, 72, 80–81
- Gordon, Robert, 73, 290n20
- Grace, Katja, 290nn23–24
- Grady v. North Carolina*, 306n50
- Graham, Howard, 151, 156–157,
 300n29, 301nn35–37, 301n44,
 312n15, 313n25
- Greely, Henry T., 195, 202, 204, 306n2,
 307nn11–12, 308nn16–17
- Greene, Jamal, 312n18
- Gregor, Mary, 308n19
- Gugliotta, Guy, 300n30
- Gunkel, David, 285n29, 296n97
- Haakonssen, Knud, 287n6, 311n2
- Halevy, Amir, 312n20
- Hal thought experiment, 9–14, 17–19,
 21–24, 28–31, 39, 45–46, 53–56, 61,
 77, 87, 90, 96, 114–118, 122, 127,
 134, 142–143, 160, 164, 184–188,
 197–199, 202, 217, 226–228, 230–
 231, 236, 239, 242, 245–247, 256,
 257–258, 263–264, 270–271, 274,
 284n28, 289n8
- Hamerooff, Stuart R., 292–293n54
- Hare, Brian, 305n28
- Harris, Allison P., 313n27
- Harrison, Harry, 287n4
- Hartmann, Thom, 300n29, 301n37
- Hassabis, Demis, 72
- Hauer, Rutger, 48
- Hauser, Marc, 173, 304nn23–24, 304n26
- Hawking, Stephen, 35, 99, 108, 240,
 275, 287n5, 294n72, 311n4
- Hayek, Friedrich, 269
- Hegel, Georg, 284n28
- Heller, Nathan, 285n29
- Helms, Jesse, 196
- Herbert, Frank, 60
Dune, 60, 62, 241
- Hermann, Peter, 281n2
- Hinton, Geoffrey, 71–72, 119–120,
 297n103
- Hitchcock, Alfred, 304n25
- Hobbes, Thomas, 26, 138, 286n36,
 299n15
- Hoffman, Reid, 286n33, 309n31
- Holland, Jesse J., 301n43
- Holland, Owen, 283n18
- Holmes, Oliver Wendell, 224
- Homenda, W., 289n14
- Horwitz, Morton, 137, 139, 299n12,
 299n17
- Hugo, Kristin, 307n10
- Humanity
 competing justifications of special
 moral status of, 8, 25–26, 168–170
 consciousness in, 26, 28
 definitions of, used to forbid research,
 198–205
 DNA percentage similarity to, 198–199
 genetic definitions of, 27, 198–199
 and human dignity, 205
 human exceptionalism, 25, 90, 114,
 117–119, 167, 177, 255
 legislative proposals to ban “blurring”
 concept of, 202–204
 as moral agents, 25–27
 personhood distinguished from, 181
 portrayal concerns about human
 likeness, 200–201

- potential-based definitions of, 201–202
procreation-based definitions of, 200
provenance of cells, 199
religious explanations of special moral status of, 27, 167–168
speciesism, 31, 201, 207, 208–210, 220–233, 249, 250, 256, 282n13
Hume, David, 237, 311n2
Hurlbut, William, 200–201
Huxley, T. H., 88–89, 168, 208, 292n47, 303n14
Hybrids, 5–6, 15–16, 31–32, 183–184, 193, 196–198, 201–203, 207–208, 214–215, 220, 226, 231, 253, 256, 270
Iida, Fumiya, 283n18
Inayatullah, Sohail, 7, 282n16
Innocent IV (pope), 131
Iwai, Katsuhito, 299n10
Jaworska, Agnieszka, 308n19
Jefferson, Thomas, 282n13
Jenkins, Ryan, 287n2
Jensen, M. C., 140, 299nn18–19
Johnson, Mark, 117, 283n18, 297n102
Johnston, Josephine, 308n18
Jones, David Albert, 308n18
Jowett, Benjamin, 303n11
Kanai, Ryota, 297n105
Kant, Immanuel, 37–39, 216, 218, 220, 302n4, 308n19
Kantian philosophy, 37–38, 115, 121, 154, 205, 217, 237, 310n44
Kass, Leon, 206, 308nn21–22, 310n42
Kavanaugh, Brett M., 314n28
Keown, Damien, 303n13
King, B. B., 112
King, Martin Luther, Jr., 273
Kleeman, Jenny, 288n25
Knight, Jack, 313n27
Korsgaard, Christine, 218
Kuhner, Timothy K., 301n52
Kuniyoshi, Yasuo, 283n18
Kurki, Visa, 299n21
Kurzweil, Ray, 65–67, 75, 82, 100, 121, 283n19, 289n9, 289n12, 298n106
Kyllo v. United States, 306n49
Lakoff, Andrew, 314nn33–34
Lakoff, George, 117, 283n18, 297n102
Lambert, Harry, 297n103
LaMDA (Language Model for Dialog Applications), 1–4, 8, 11, 71, 85, 95, 242, 274, 281n4, 283n22
Language
Aristotle's theory of, 117, 166–167, 244
chatbot's fluency in, impact of, 3–4, 12, 25, 110–118, 167, 242–243, 272
“computational shallowness” of, 12, 25, 95, 119, 245, 272
human exceptionalism, last citadel of, 25–26, 113–114, 118–120, 167, 218
learning through embodied intelligence, 10
quality distinguishing humans from other animals, 1, 21, 24–26, 87, 115, 117–118, 166–168, 171–172, 178, 202, 216–219, 230, 242, 249
and sentience, implications of, 30, 111, 118
Language, last citadel of, 25, 118–120 and human exceptionalism, 25, 113, 118–120, 167, 218
Large language models, 1, 5, 10–11, 71, 88, 95, 111, 115, 119, 240–241, 243–245, 283n22
Lau, Hakwan, 294n67
Lawrence v. Texas, 306n55
Leahy, Patrick, 284n27
LeCun, Yann, 119, 297n103
Lee, James J., 304n25
Lee, Raymond, 283n22

- Le Guin, Ursula K., 42, 287n18
- Lemley, Mark, 314n30
- Lemoine, Blake, 1–5, 11, 34, 71, 85, 95, 108, 111, 121–122, 243, 261, 281n5
- Lenat, Doug, 73
- Lenharo, Mariana, 293n60, 297n105
- Lessig, Larry, 300n28
- Lin, Patrick, 287n2
- Litt, Abninder, 293n55
- Lizza, John P., 312n20
- LLMs. *See* Large language models
- Lohr, Steve, 285n29
- Long, Robert, 293n64, 294n68
- Lonsdorf, Elizabeth V., 304n28
- Louden, Robert, 302n4
- Lovecraft, H. P., 240
- Loving v. Virginia*, 209, 271, 309nn28–29
- MacKellar, Calum, 308n18
- Macukow, Bohdan, 289n14
- Magrath, C. Peter, 300n29
- Maitland, Frederic, 136, 142
- Malcolm X, 273
- Manjoo, Farhad, 295n93
- Marr, Bernard, 281n8
- McCain, John, 144–145
- McCarthy, John, 289n7
- Mcgraw, Peter, 304n17
- McMahan, Jeff, 309n26
- Meckling, W. H., 140, 299nn18–19
- Mehendale, H. M., 307n6
- Melis, Alicia P., 305n28
- Melone, Matthew A., 300n52
- Meredith, Richard, 288n31
- Metz, Cade, 290n17, 290n22, 294n76
- Miceli, Thomas J., 313n27
- Michel, Matthias, 294n67
- Midjourney, 118
- Minsky, Marvin, 64, 289n7, 295n84
- Miracchi Titus, Lisa, 283n18
- Moral philosophy
- contradiction and tension as productive parts of, 237–238
 - contributions to personhood debate, 7, 18–21, 35–39, 163–165, 185, 205–234, 284n28
 - critique of universal moral theories, 236–237
 - empathy's connection to, 35–41, 238–239
- Morgenbesser, Sidney, 91
- Morris, Thomas D., 299n21
- Musk, Elon, 102, 294nn75–76, 294n78
- Naruto v. Slater*, 183, 184, 188, 305n43, 306n45
- Nat'l Mut. Ins. Co. of D.C. v. Tidewater Transfer Co.*, 301n46
- Naturalistic fallacy, 29
- New Yorker*, 25, 285n29, 311n6
- New York Times*, 4, 13, 51, 145, 281n9, 283nn23–24, 285n29, 288n26, 290n17, 290n22, 292n52, 294n66, 294n76, 295n93, 300n30, 307n5, 311n7
- Nexus theory, 140
- Ng, Andrew, 102, 294n77
- Nguyen, Phuoc, 285n29
- Nim the chimpanzee, 305n30
- Nodelman, Uri, 308n19, 310n49
- Nonhuman Rights Project, 163, 165, 302nn1–2
- Norvig, Peter, 12, 71, 80, 109, 243, 283n21, 289n15, 295n84, 295n91, 311n9
- Obergefell v. Hodges*, 158, 301nn50–51, 306nn56–57
- Obeyesekere, Gananath, 303n13
- Oliver, John, 44
- Palacios-González, César, 308n18, 308n25, 309n30
- Palila v. Hawaii Dep't of Land & Nat. Res.*, 305n41
- Parker, Dorothy, 174

- Pathetic fallacy, 34
- Paulsen, Michael Stokes, 313n26
- Penrose, Roger, 293n54
- People ex rel. Nonhuman Rts. Project v. Lavery*, 303n7, 305n39
- Pepperberg, Irene M., 303n17
- Personhood
- and anthropomorphism, 6, 29, 33–35
 - corporate, 131–161
 - critiques of AI personhood debate, 7–8
 - definitions of, competing, 27–28, 37–42, 49–57, 65, 74, 83–85, 108–110, 114–116, 120–130
 - designing entities around definitions of, 55–56
 - fetal, 19, 189, 238, 259–263, 269, 313n24, 313n26
 - line between persons and nonpersons, 5–9, 15–16, 18–32, 37–42, 49–57, 65, 74, 83–85, 108–110, 114–116, 120–130, 241, 257
 - nonhuman animals and, 163–193
 - science fiction explorations of, 31, 41–57
 - transgenic entities, chimeras, and hybrids, and, 195–234
- Petrovski, Henry, 286n32
- Pfeifer, Rolf, 283n18
- Phillips, Michael J., 299n16
- Piaget, Jean, 296n101
- Piattelli-Palmarini, Massimo, 296n101
- Pichai, Sundar, 281n10
- Pierce, Jessica, 303n17
- Pinker, Stephen, 304n25
- Piotrowska, Monika, 307n7
- Planned Parenthood v. Casey*, 314n28
- Platt, Charles, 77
- Platt's Law, 77, 79
- Pope, Alexander, 245, 312n11
- Poynder, John, 300n25
- Psychological priming, 49–58, 262–274, 288n27
- Putnam, Hilary, 293n65
- Quinn, Warren, 308n19
- Rabin, Roni Caryn, 283n24
- Race, 32, 66, 148, 256, 269, 312n18
- racism, 31, 112, 149, 153, 189, 207–210, 224–225, 233, 255–256, 282n13
 - as a social construct, 121, 170, 209, 217, 220, 224, 227, 233, 268
 - white supremacy, 7, 43, 149, 209, 215
- Rand, Ayn, 7, 269
- Ray, C. Claiborne, 283n23
- Real entity theory, 126, 134, 136–137, 142–143, 180, 227, 247, 299n15
- Regan, Tom, 308n19
- Rehnquist, William, 145–147, 152, 157–159, 181, 247–248, 251, 300nn26–28, 301n45, 312nn13–14
- Rescorla, Michael, 293n57
- Robbe-Grillet, Alain, 34
- Robert, Jason Scott, 310nn37–38
- Roberts, John, 306n57, 314n28
- Rochester, Nathaniel, 289n7
- Roden, Gregory J., 313n24
- Roe v. Wade*, 155, 188, 200, 259, 266, 314n28
- Rogers, G. A. J., 286n36, 299n15
- Romilly, Samuel, 40
- Roose, Kevin, 4, 281n9, 282n11, 295n93, 311n7
- Rosati, Alexandra G., 305n28
- Rosen, Jeffrey, 281n6, 311n52, 312n19
- Rosin, Michael L., 313n24
- Roth, Philip, 43–44
- Rudin, Cynthia, 298n108
- Ruskin, John, 34, 287n3
- Russell, Stuart, 12, 80, 109, 243, 283n21, 295n84, 295n91, 311n9
- Ryerson, James, 292n52
- Saeed, K., 289n14
- Samsa, Gregor, 207, 227, 256–257
- San Mateo v. Southern Pacific Railroad Co.*, 299n17, 313n25

- Santa Clara County v. Southern Pacific Railroad*, 151–154, 156
- Scalia, Antonin, 156, 186, 300n28, 301n43, 306n52
- Schacter, Daniel L., 288n27
- Schafer, Burkhard, 133
- Schmidhuber, Jürgen, 79, 82, 289n14
- Schuhmann, Karl, 286n36, 299n15
- Science fiction, 6, 18–19, 41–57
explorations of artificial personhood in, 18, 31, 41–57
- Scott, Ridley, 35, 41, 46–47, 49, 52, 239, 262
- Seager, William, 289n5
- Searle, John, 5, 84–85, 87, 88–94, 96–98, 110–115, 120–121, 142, 242–243, 245–246, 262, 272, 274, 292n43, 292nn50–51, 292n53
- Sen, Maya, 313n27
- Shagrir, Oron, 293n65
- Shakespeare, William, 39, 297n104
- Shankar, K., 307n6
- Shannon, Claude E., 289n7
- Shapiro, Lawrence, 297n101
- Shepherd, Joel, 295n88, 310n45
- Shewmon, D. Alan, 312n20
- Shreve, Jamie, 307n5, 307nn8–9
- Shukla, Samarth, 296n100
- Siegel, Andrew, 308n24
- Silver, David, 291n38
- Simon, Herbert, 64, 289n8
- Singer, Peter, 206, 218, 237, 249, 256, 311n1
- Singularity, the, 65–67, 76, 78, 99, 100–101, 108, 175, 239, 289n9
- Sisk, Gregory C., 313n27
- Skinner, B. F., 26, 42, 46, 86, 91, 242, 286n35, 287n23, 292nn45–46, 311n8
- Skynet, 7, 19, 98
- Slaughter-House Cases*, 300n32
- Slavery, 16, 19, 64, 147, 148, 213, 287n15, 299n21
- abolition of, 39
- British Society for the Abolition of Slavery, 6
- and enslaved persons, 7, 20, 38, 40, 99, 143, 150, 185–186, 229, 248, 264, 284n28
- slave revolts, 45
- Smith, Adam, 16, 30, 35–36, 38, 40–41, 45, 49, 53–55, 109, 121, 218, 237–239, 251, 287nn16–17, 288n29, 312n17
- empathy as basis for morality, 37–39
- Theory of the Moral Sentiments*, 30, 35–38, 40–41, 218, 238, 287n6, 311n2
- Smith, Lamar, 284n27
- Smith, N. R., 306n45
- Smyth, Nicholas, 310n49
- Solaiman, S. M., 285n29
- Solum, Lawrence, 129, 284n29, 298n112
- Somers, James, 311n6
- Souls, 27, 112, 168
human, 104, 168
- soullessness, 33, 136, 145, 168, 260, 268
- Spaulding, Shannon, 297n101
- Spinoza, Baruch, 39
- Stable Diffusion, 118
- Standish, Paul, 297n104
- Steels, Luc, 283n18
- Stein, F. N., 15–17, 20, 234
- Stepney, S., 296n99
- Stevens, John Paul, 145, 158–159, 264, 268, 300n23, 312n16, 314n31
- Stone, Christopher, 7, 282n16, 305n44
- Story, Joseph, 131–133, 142, 146, 247
- Stowe, Harriet Beecher, 40
Uncle Tom's Cabin, 40, 160
- Streiffner, Robert, 308n20, 309n32, 310nn39–40, 310n44
- Stuart Little*, 200
- Subbaraman, Nidhi, 307n4

- Tabarrok, Alex, 297n103
Taddeo, Sarah, 307n3
Taney, Roger, 185
Tannenbaum, Julie, 308n19
Tarantola, Andrew, 287n24
Taylor, Robert M., 312n21
Technological change
 difficulty in predicting speed of, 68–70, 72–83, 101
 linear as opposed to exponential, 66–68, 100–101
Tennessee Valley Authority v. Hiram Hill et al., 295n87
Terminator, The (film), 19, 52, 269
Terrace, Herbert S., 305n30
Tests of sentience, 28
 and Chinese Room (*see* Chinese Room)
 and Turing Test, 21, 30, 35, 45–46, 95
 (*see also* Turing Test)
 and Voight-Kampff Test, 30, 47, 49
 (*see also* Voight-Kampff Test)
Teubner, Gunther, 286n29
Thaler v. Vidal, 296n94
Theunissen, L. Nandi, 308n19
Thompson, Evan, 296n101
Thurber, James, 174
Thurlow, Edward, 145
Tiku, Natasha, 281n1, 281n4, 281n7
Tilden, Mark, 33
Timmermann, Jens, 308n19
To Kill a Mockingbird, 36
Tomasello, Michael, 305n28
Tommy the chimpanzee, 163–165, 175–178, 180–188, 192–193, 225–226, 228, 231, 250–253, 302n1, 303n7, 313n24
Tomson, Keith, 303n14
Tononi, Giulio, 95–96, 293n58, 293n63
Transgenic species, 8–9, 14, 18–22, 31, 134, 163–165, 170, 183–186, 195–198, 203–205, 207–208, 212–213, 220, 226, 235–236, 253, 255, 259, 261, 270, 282n13, 285n29, 307n6
Chimpy thought experiment (*see* Chimpy thought experiment)
OncoMouse, 17, 21, 196, 203, 271
Treguer, Pascal, 304n25
Tricker, Robert Ian, 299n18
Trump, Donald, 190, 235
Truog, Robert D., 312n21
Trustees of Dartmouth College v. Woodward, 298n2, 312n12
Tulving, Endel, 288n27
Turing, Alan, 25, 46, 56, 63, 83–86, 95, 109, 111, 115, 118, 287n22, 291nn41–42, 292n45
Turing Test, 10–13, 46, 63–65, 75, 83–86, 108–121, 242–244, 283n19, 286n33, 291n42
advantages and limitations of, 108–109, 110–121, 243–244, 286n33
in the age of chatbots, 110–121
Tushnet, Mark, 299n21
Twain, Mark, 299n20, 304n19
2001: A Space Odyssey, 19
Ulam, Stanisław, 66, 289nn12–13
Updike, John, 43–44
van Dijk, Jelle, 282n13
van Gogh, Vincent, 112
VanRullen, Rufin, 297n105
Verma, Pranshu, 282n10
Vincent, James, 294n78, 294n80, 298n5
Vinge, Vernor, 65–77, 100–101
Virtual assistants
 Alexa, 21
 Siri, 1, 5, 14, 18, 21, 51, 63, 83, 94, 239
Vitiello, G., 292n54
Voight-Kampff Test, 43–46, 55, 57, 98, 109, 114, 120, 160, 218, 239, 242
encounters with AI as version of, 57, 109, 114, 120, 160

- Voight-Kampff Test (cont.)
 humanity's ability to pass our own,
 160
 irony of, 44, 57
 metaphor for defining humanity, 57,
 109, 114, 120, 160
 test of empathy for nonhuman
 animals, 43–44
 test to identify, and kill, androids, 44
 von Bismarck, Otto, 189
 von Gierke, Otto, 136, 142
 von Neumann, John, 66–67
 von Wright, G. H., 297n104
- Waite, Morrison, 151–152, 300n29
 Waldau, Paul, 303n13
Wall Street Journal, 14
 Walsh, R., 296n99
 Warner, Joel, 304n17
 Warren, Earl, 189, 209
 Washington, Booker T., 273
Washington Post, 1–2, 13, 281nn1–2,
 282n10, 287n1, 298n109
 Watson, John B., 292n45
 “We, Robot” Conference, 278
 Webster, Daniel, 287n15
 Weissman, Irving, 195, 199, 201–202,
 204, 271
 Weld, Theodore Dwight, 39
 Whang, Oliver, 294n66, 294n69
 Wheeler, Lydia, 313n22
Wheeling Steel v. Glander, 158, 301n47
 White, Stanford, 304n25
 Wilberforce, Samuel, 88, 168, 292n47,
 303n14
 Wilberforce, William, 40
 Williams, Bernard, 223–224, 303n6,
 310nn48–49
 Williams, Chris, 294n77
 Winfield, A. F., 296n99
 Winkler, Adam, 153, 300n29, 301n37,
 301n39, 301n42
 Wise, Steven, 163, 165
- Wittes, Benjamin, 281n6, 311n52,
 312n19
 Wittgenstein, Ludwig, 26, 119, 286n36,
 297n104
 Wolfe, Gene, 287n4
 Wolfram, Stephen, 12, 25, 245, 283n20,
 286n34, 312n10
 Wyllie, David, 286n30
- Youn, Monica, 301n52
 Yudkowsky, Eliezer, 66, 75, 82, 102,
 104–105, 107–108, 116, 275,
 289n11, 295n85, 295nn89–90
- Zuckerberg, Mark, 102, 294n76
 Zuse, Konrad, 79