CONTENTS

Preface to the Paperback Edition xi

1 Introduction	1
The Dawn of AI as a Consumer Product	3
AI Shakes Up Entertainment	7
Predictive AI: An Extraordinary Claim That Requires Extraordinary Evidence	9
Painting AI with a Single Brush Is Tempting but Flawed	12
A Series of Curious Circumstances Led to This Book	18
The AI Hype Vortex	21
What Is AI Snake Oil?	26
Who This Book Is For	34
2 How Predictive Al Goes Wrong	36
Predictive AI Makes Life-Altering Decisions	38
A Good Prediction Is Not a Good Decision	43
Opaque AI Incentivizes Gaming	46
Overautomation	48
Predictions about the Wrong People	51

For general queries, contact info@press.princeton.edu.

viii CONTENTS

	Predictive AI Exacerbates Existing Inequalities	53
	A World without Prediction	56
	Concluding Thoughts	58
3	Why Can't Al Predict the Future?	60
	A Brief History of Predicting the Future	
	Using Computers	62
	Getting Specific	67
	The Fragile Families Challenge	70
	Why Did the Fragile Families Challenge End	
	in Disappointment?	73
	Predictions in Criminal Justice	78
	Failure Is Hard. What about Success?	81
	The Meme Lottery	86
	From Individuals to Aggregates	90
	Recap: Reasons for Limits to Prediction	97
4	The Long Road to Generative Al	99
	Generative AI Is Built on a Long Series	
	of Innovations Dating Back Eighty Years	105
	Failure and Revival	107
	Training Machines to "See"	111
	The Technical and Cultural Significance of ImageNet	114
	Classifying and Generating Images	118
	Generative AI Appropriates Creative Labor	122
	AI for Image Classification Can Quickly	
	Become AI for Surveillance	127

CONTENTS ix

	From Images to Text	129
	From Models to Chatbots	133
	Automating Bullshit	139
	Deepfakes, Fraud, and Other Malicious Uses	142
	The Cost of Improvement	143
	Taking Stock	146
5	Is Advanced AI an Existential Threat?	150
	What Do the Experts Think?	151
	The Ladder of Generality	156
	What's Next on the Ladder?	162
	Accelerating Progress?	165
	Rogue AI?	168
	A Global Ban on Powerful AI?	172
	A Better Approach: Defending against	
	Specific Threats	174
	Concluding Thoughts	177
6	Why Can't Al Fix Social Media?	179
	When Everything Is Taken Out of Context	183
	Cultural Incompetence	188
	AI Excels at Predicting the Past	194
	When AI Goes Up against Human Ingenuity	198
	A Matter of Life and Death	201
	Now Add Regulation into the Mix	205
	The Hard Part Is Drawing the Line	209

X CONTENTS

	Recap: Seven Shortcomings of AI for	
	Content Moderation	216
	A Problem of Their Own Making	218
	The Future of Content Moderation	223
7	Why Do Myths about AI Persist?	227
	AI Hype Is Different from Previous Technology Hype	231
	The AI Community Has a Culture and History of Hype	235
	Companies Have Few Incentives for Transparency	239
	The Reproducibility Crisis in AI Research	241
	News Media Misleads the Public	247
	Public Figures Spread AI Hype	251
	Cognitive Biases Lead Us Astray	255
8	Where Do We Go from Here?	258
	AI Snake Oil Is Appealing to Broken Institutions	261
	Embracing Randomness	265
	Regulation: Cutting through the False Dichotomy	268
	Limitations of Regulation	274
	AI and the Future of Work	276
	Growing Up with AI in Kai's World	281
	Growing Up with AI in Maya's World	285

Epilogue to the Paperback Edition 291

Acknowledgments 301

References 303

Index 343
For general queries, contact info@press.princeton.edu.

© Copyright Princeton University Press. No part of this book may be
distributed, posted, or reproduced in any form by digital or mechanical
means without prior written permission of the publisher.

C	h	ล	n	t	ρ	r	1
U	Ш	u	μ	ι	U	ı	

IMAGINE AN ALTERNATE universe in which people don't have words for different forms of transportation—only the collective noun "vehicle." They use that word to refer to cars, buses, bikes, spacecraft, and all other ways of getting from place A to place B. Conversations in this world are confusing. There are furious debates about whether or not vehicles are environmentally friendly, even though no one realizes that one side of the debate is talking about bikes and the other side is talking about trucks. There is a breakthrough in rocketry, but the media focuses on how vehicles have gotten faster—so people call their car dealer (oops, vehicle dealer) to ask when faster models will be available. Meanwhile, fraudsters have capitalized on the fact that consumers don't know what to believe when it comes to vehicle technology, so scams are rampant in the vehicle sector.

Now replace the word "vehicle" with "artificial intelligence," and we have a pretty good description of the world we live in.

Artificial intelligence, AI for short, is an umbrella term for a set of loosely related technologies. ChatGPT has little in common with, say, software that banks use to evaluate loan applicants. Both are referred to as AI, but in all the ways that matter—how

they work, what they're used for and by whom, and how they fail—they couldn't be more different.

Chatbots, as well as image generators like Dall-E, Stable Diffusion, and Midjourney, fall under the banner of what's called generative AI. Generative AI can generate many types of content in seconds: chatbots generate often-realistic answers to human prompts, and image generators produce photorealistic images matching almost any description, say "a cow in a kitchen wearing a pink sweater." Other apps can generate speech or even music.

Generative AI technology has been rapidly advancing, its progress genuine and remarkable. But as a product, it is still immature, unreliable, and prone to misuse. At the same time, its popularization has been accompanied by hype, fear, and misinformation.

In contrast to generative AI is predictive AI, which makes predictions about the future in order to guide decision-making in the present. In policing, AI might predict "How many crimes will occur tomorrow in this area?" In inventory management, "How likely is this piece of machinery to fail in the next month?" In hiring, "How well will this candidate perform if hired for this job?"

Predictive AI is currently used by both companies and governments, but that doesn't mean it works. It's hard to predict the future, and AI doesn't change this fact. Sure, AI can be used to pore over data to identify broad statistical patterns—for instance, people who have jobs are more likely to pay back loans—and that can be useful. The problem is that predictive AI is often sold as far more than that, and it is used to make decisions about people's lives and careers. It is in this arena that most AI snake oil is concentrated.

AI snake oil is AI that does not and cannot work as advertised. Since AI refers to a vast array of technologies and applications,

most people cannot yet fluently distinguish which types of AI are actually capable of functioning as promised and which types are simply snake oil. This is a major societal problem: we need to be able to separate the wheat from the chaff if we are to make full use of what AI has to offer while protecting ourselves from its possible harms, harms which in many cases are already occurring.

This book is a guide to identifying AI snake oil and AI hype. In it, we'll give you essential vocabulary to tease apart generative AI, predictive AI, and other types of AI. We'll share commonsense ways of assessing whether or not a purported advance is plausible. This will make you read news about AI much more skeptically and with an eye toward details that often get buried. A deeper understanding of AI will both satisfy your scientific curiosity and translate into practical ideas on how to use—and when not to use—AI in your life and career. And we will make the argument that predictive AI not only does not work today but will likely never work, because of the inherent difficulties in predicting human behavior. Finally, we hope that this book will get you thinking about your own responsibilities—and opportunities for change—with respect to the harmful implications of these tools.

The Dawn of AI as a Consumer Product

Generative AI burst into public view with the release of Chat-GPT in November 2022. Ironically, the developer OpenAI released it without any fanfare and internally thought of it as a "research preview." But it went viral overnight as people shared amusing and amazing examples of its capabilities. For example, software developer Thomas Ptacek asked it to write a biblical verse in the style of the King James Bible, explaining how to

remove a peanut butter sandwich from a VCR.² This is what ChatGPT spat out:

And it came to pass that a man was troubled by a peanut butter sandwich, for it had been placed within his VCR, and he knew not how to remove it.

And he cried out to the Lord, saying, "Oh, Lord, how can I remove this sandwich from my VCR, for it is stuck fast and will not budge?"

And the Lord spoke unto him, saying, "Fear not, my child, for I shall guide thy hand and show thee the way. Take thy butter knife, and carefully insert it between the sandwich and the VCR, and gently pry them apart. And with patience and perseverance, the sandwich shall be removed, and thy VCR shall be saved."

And the man did as the Lord commanded, and lo and behold, the sandwich was removed from the VCR, and the man was saved.

And the Lord said, "Verily I say unto thee, seek not to put thy peanut butter sandwiches in thy VCR, for it is not a suitable place for such things. Rather, keep thy sandwiches in thy refrigerator or on thy plate, where they belong."

And the man heeded the Lord's words, and from that day forth, he kept his sandwiches in their proper place, and was saved from trouble and woe. Amen.

Two months later, the app reportedly had over one hundred million users.³ OpenAI was caught off guard by the explosion of interest and hadn't even procured enough computing power to handle the traffic that it generated.

Computer programmers soon started using it, because it turned out that ChatGPT was pretty good at generating snippets of programming code given only a description of what the

code should do. In fact, programmers had been using a previous product called GitHub Copilot that used similar technology, but their use of AI accelerated after ChatGPT's release. It decreased the time it took to create apps. Even nonprogrammers could now build some simple apps.

Microsoft soon licensed the technology from OpenAI, and its Bing search engine rolled out a chatbot that could answer questions based on search results. Google had built its own chatbot in 2021 but had not yet released it or integrated it into its products. Bing's move was seen as an existential threat to Google, and Google hurriedly announced its own search chatbot called Bard (later renamed Gemini).

That's when things started to go wrong. In the promotional video for Bard, the bot said that the James Webb Space Telescope took the first picture of a planet outside the solar system. An astrophysicist pointed out that this was wrong. Apparently Google couldn't get even a cherry-picked example right. Its market value instantly took a hundred-billion-dollar dip. That's because investors were spooked by the prospect of a search engine that would get much worse at answering simple factual queries if Google were to integrate Bard into search, as it had promised.

Google's embarrassment, while expensive, was only a ripple that portended the wave of problems that arose from chatbots' difficulties with factual information. Their weakness is a consequence of the way they are built. They learn statistical patterns from their training data—which comes largely from the web—and then generate remixed text based on those patterns. But they don't necessarily remember what's in their training data. We'll dive into this in chapter 4.

Misuse of the technology is rampant. News websites have been caught publishing error-filled AI-generated stories on

important topics such as financial advice, and then refusing to stop using the technology even after the errors came to light.⁷ Amazon is overrun with AI-generated books, including a few mushroom foraging guides, where errors can be fatal if a reader trusts the book.⁸

It's easy to look at all the flaws and misuses of chatbots and conclude that the world has gone mad for being so gaga about a technology that is so failure prone. But that conclusion would be too simplistic.

We think most knowledge industries can benefit from chatbots in some way. We use them ourselves for research assistance, for tasks ranging from mundane ones such as formatting citations correctly, to things we wouldn't otherwise be able to do such as understanding a jargon-filled paper in a research area we aren't familiar with.

The catch is that it takes effort and practice to use chatbots while avoiding their ever-present pitfalls. But *inappropriate* uses are much easier, because someone trying to make a quick buck, say by selling an AI-generated book, doesn't often care if the contents are garbage. That's what makes chatbots so conducive to misuse.

There are thornier questions about power. Suppose web search companies replace their traditional list of ten links with AI-generated ready answers. Even assuming that accuracy problems are fixed, the result is basically a machine for rewriting content found on other websites and passing it off as original, without having to send traffic or revenue to those websites. If search engines simply presented others' content as their own, they would run afoul of copyright law. But AI-generated answers seem to skirt this issue, although there are many lawsuits seeking to change this as of 2024.⁹

AI Shakes Up Entertainment

Another generative AI technology that has captivated people is text-to-image generation. In mid-2023, it was estimated that over a billion images had been created using Dall-E 2 by OpenAI, Firefly by Adobe, and Midjourney (by a company of the same name). Another widely used image generator is Stable Diffusion by Stability AI, which is openly available, meaning that anyone can modify it to their liking. Stable Diffusion—based tools have been *downloaded* over two hundred million times. Since users run it on their own devices, there is no central tally of how many images have been generated using it, but it is likely to be several billion.

Image generators have enabled a deluge of entertainment. ¹¹ Unlike traditional entertainment, these images are endlessly customizable to each user's interests. Some people delight in fantastic landscapes or cityscapes. Others enjoy images of historical figures in modern situations, or famous people doing things they wouldn't normally do, such as the Pope wearing a puffer jacket, dubbed "Balenciaga Pope." Fake trailers for various movies such as *Star Wars* in the highly recognizable style of Wes Anderson—symmetrical framing, pastel colors, whimsical sets—have proven popular.

It's not only hobbyists who are excited about image generators: entertainment apps are big business. Video game companies have created in-game characters that players can have a natural conversation with.¹² Many photo editing apps now have generative AI functionality. So, for example, you can ask such an app to add balloons to a picture of a birthday party.

AI was a major point of contention in the 2023 Hollywood strikes. 13 Actors worried that studios would be able to use

existing footage of them to train AI tools capable of generating new videos based on a script—videos that looked like they featured the real actors whose images and videos the AI tools were trained on. In other words, studios would be able to capitalize on actors' likenesses and past labor in perpetuity, but without compensation.

While the strikes have ended, the underlying tensions between labor and capital are sure to resurface, especially as the technology advances. ¹⁴ Many companies are working on text-to-video generators, while others are working on automating script writing. The end result might not be as artistically complex or valuable, but that might not matter to studios looking to crank out a summer blockbuster.

In the long run, we think that a combination of technology and law can alleviate most of the problems we've described, as well as amplify the benefits. For example, there are many promising technical ideas to make chatbots less likely to fabricate information, while regulation can curb intentional misuses. But in the short term, adjusting to a world with generative AI is proving to be painful, as these tools are highly capable but unreliable. It's as if everyone in the world has been given the equivalent of a free buzzsaw.

It will take work to integrate AI appropriately into our lives. A good example is what's happening in schools and colleges, given that AI can generate essays and pass college exams. Let's be clear—AI is no threat to education, any more than the introduction of the calculator was. With the right oversight, it can be a valuable learning tool. But to get there, teachers will have to overhaul their curricula, their teaching strategies, and their exams. At a well-funded institution such as Princeton, where we teach, this is an opportunity rather than a challenge. In fact, we encourage our students to use AI. But many others have

been left scrambling as ChatGPT suddenly put a potential cheating tool in the hands of millions of students.

Will society be left perpetually reacting to new developments in generative AI? Or do we have the collective will to make structural changes that would allow us to spread out the highly uneven benefits and costs of new innovations, whatever they may be?

Predictive AI: An Extraordinary Claim That Requires Extraordinary Evidence

Generative AI creates many social costs and risks, especially in the short term. But we're cautiously optimistic about the potential of this type of AI to make people's lives better in the long run. Predictive AI is a different story.

In the last few years, applications of predictive AI to predict social outcomes have proliferated. Developers of these applications claim to be able to predict future outcomes about people, such as whether a defendant would go on to commit a future crime or whether a job applicant would do well at a job. In contrast to generative AI, predictive AI often does not work at all.¹⁶

People in the United States over the age of sixty-five are eligible to enroll in Medicare, a state-subsidized health insurance plan. To cut costs, Medicare providers have started using AI to predict how much time a patient will need to spend in a hospital. These estimates are often incorrect. In one case, an eighty-five-year-old was evaluated as being ready to leave in seventeen days. But when the seventeen days passed, she was still in severe pain, and couldn't even push a walker without help. Still, based on the AI assessment, her insurance payments stopped. In cases like this, AI technology is often deployed with sensible

intentions. For example, without predictive AI, nursing homes would be logically incentivized to house patients forever. But in many cases, the goals of the system as well as how it's deployed change over time. One can easily imagine how Medicare providers' use of AI may have started as a way to create a modicum of accountability for nursing homes, but then morphed into a way to squeeze pennies out of the system regardless of the human cost.

Similar stories are prevalent across domains. In hiring, many AI companies claim to be able to judge how warm, open, or kind someone is based on their body language, speech patterns, and other superficial features in a thirty-second video clip. Does this really work? And do these judgments actually predict job performance? Unfortunately, the companies making these claims have failed to release any verifiable evidence that their products are effective. And we have lots of evidence to the contrary, showing that it is extremely hard to predict individuals' life outcomes, as we'll see in chapter 3.

In 2013, Allstate, an insurance company, wanted to use predictive AI to determine insurance rates in the U.S. state of Maryland—so that the company could make more money without losing too many customers. It resulted in a "suckers list"—a list of people whose insurance rates increased dramatically compared to their earlier rates. Seniors over the age of sixty-two were drastically overrepresented in this list, an example of automated discrimination. It is possible that seniors are less likely to shop around for better prices and that AI picked up on that pattern in the data. The new pricing would likely increase revenue for the insurance company, yet it is morally reprehensible. While Maryland refused Allstate's proposal to use this AI tool on the grounds that it was

discriminatory, the company does use it in at least ten other U.S. states.*

If individuals object to AI in hiring, they can simply choose not to apply for jobs that engage AI to judge résumés. When predictive AI is used by governments, however, individuals have no choice but to comply. (That said, similar concerns also arise if many companies were to use the same AI to decide who to hire.) Many jurisdictions across the world use criminal risk prediction tools to decide whether defendants arrested for a crime should be released before their trial. Various biases of these systems have been documented: racial bias, gender bias, and ageism. But there's an even deeper problem: evidence suggests that these tools are only slightly more accurate than randomly guessing whether or not a defendant is "risky."

One reason for the low accuracy of these tools could be that data about certain important factors is not available. Consider three defendants who are identical in terms of the features that might be used by predictive AI to judge them: age, the number of past offenses, and the number of family members with criminal histories. These three defendants would be assigned the same risk score. However, in this example, one defendant is deeply remorseful, another has been wrongly arrested by the police, and the third is itching to finish the job. There is no good way for an AI tool to take these differences into account.

Another downside of predictive AI is that decision subjects have strong incentives to game the system. For example, AI was used to estimate how long the recipient of a kidney transplant

^{*} Many of the examples in this book, like this one, are from the United States, simply because that is where we are based. However, the lessons we draw from these examples are intended to be broadly applicable.

would live after their transplant.¹⁹ The logic was that people who had the longest to live after a transplant should be prioritized to receive kidneys. But the use of this prediction system would *disincentivize* patients with kidney issues to take care of their kidney function. That's because if their kidneys failed at a younger age, they would be more likely to get a transplant! Fortunately, the development of this system involved a deliberative process with participation by patients, doctors, and other stakeholders. So, the incentive misalignment was recognized and the use of predictive AI for kidney transplant matching was abandoned.

We'll see many more failures of predictive AI in chapters 2 and 3. Are things likely to improve over time? Unfortunately, we don't think so. Many of its flaws are inherent. For example, predictive AI is attractive because automation makes decision-making more efficient, but efficiency is exactly what results in a lack of accountability. We should be wary of predictive AI companies' claims unless they are accompanied by strong evidence.

Painting AI with a Single Brush Is Tempting but Flawed

Generative and predictive AI are two of the main types of AI. How many other types of AI are there? There is no way to answer that question, since there is no consensus about what is and isn't AI.

Here are three questions about how a computer system performs a task that may help us determine whether the label AI is appropriate. Each of these questions captures something about what we mean by AI, but none is a complete definition. First, does the task require creative effort or training for a human to perform? If yes, and the computer can perform it, it might be AI. This would explain why image generation, for example, qualifies

as AI. To produce an image, humans need a certain amount of skill and practice, perhaps in the creative arts or in graphic design. But even *recognizing* what's in an image, say a cat or a teapot—a task that is trivial and automatic for humans—proved daunting to automate until the 2010s, yet object recognition has generally been labeled AI. Clearly, comparison to human intelligence is not the only relevant criterion.

Second, we can ask: Was the behavior of the system directly specified in code by the developer, or did it indirectly emerge, say by learning from examples or searching through a database? If the system's behavior emerged indirectly, it might qualify as AI. Learning from examples is called machine learning, which is a form of AI. This criterion helps explain why an insurance pricing formula, for example, might be considered AI if it was developed by having the computer analyze past claims data, but not if it was a direct result of an expert's knowledge, even if the actual rule was identical in both cases. Still, many manually programmed systems are nonetheless considered AI, such as some robot vacuum cleaners that avoid obstacles and walls.

A third criterion is whether the system makes decisions more or less autonomously and possesses some degree of flexibility and adaptability to the environment. If the answer is yes, the system might be considered AI. Autonomous driving is a good example—it is considered AI. But like the previous criteria, this criterion alone can't be considered a complete definition—we wouldn't call a traditional thermostat AI, one that contains no electronics. Its behavior rather arises from the simple principle of a metal expanding or contracting in response to changes in temperature and turning the flow of current on or off.

In the end, whether an application gets labeled AI is heavily influenced by historical usage, marketing, and other factors. We won't fret about the fact that there's no consistent definition.

That might seem surprising for a book about AI. But recall our overarching message: there's almost nothing one can say in one breath that applies to all types of AI. Most of our discussion in the book will be about specific types of AI, and as long as each type is clearly defined, we'll be on the same page.

There's a humorous AI definition that's worth mentioning, because it reveals an important point: "AI is whatever hasn't been done yet." In other words, once an application starts working reliably, it fades into the background and people take it for granted, so it's no longer thought of as AI. There are many examples: Robot vacuum cleaners like the Roomba. Autopilot in planes. Autocomplete on our phones. Handwriting recognition. Speech recognition. Spam filtering. Spell-check. Yes, there was a time when spell-check was considered a hard problem!

We think these tools are all wonderful. They quietly make our lives better. These are the kinds of AI we want more of. This book is about the types of AI that are problematic in some way, because you wouldn't want to read three hundred pages on the virtues of spell-check. But it's important to recognize that not all AI is problematic—far from it.

Some new AI technologies will hopefully one day come to be seen as mundane. Today, self-driving cars often make the news for accidents and fatalities. ²⁰ But safe automated driving is ultimately a solvable problem, although one whose difficulty has repeatedly been underestimated. The bigger challenge for society might be the massive labor displacement that the technology will cause if it becomes widespread—millions of people drive trucks, taxis, or rideshare vehicles. Still, if the safety problem is solved and the necessary social and political adjustments are made, we may one day take self-driving cars for granted, like we do elevators today.

However, we think other types of AI, notably predictive AI, are unlikely to become normalized. Accurately predicting people's social behavior is not a solvable technology problem, and determining people's life chances on the basis of inherently faulty predictions will always be morally problematic.

For a more in-depth case study of why we must avoid sweeping generalizations about AI, consider facial recognition, an AI technology that has civil liberties advocates concerned. It has led to many false arrests in the United States—six, as we write this—all Black people. Should the use of facial recognition by police be discontinued because it is error prone and misidentifies Black people more often?

One fact that's easy to miss in this debate is that all the false arrests involved a cascading set of police failures, most of them human errors rather than technological. Robert Williams was arrested for shoplifting in part based on the testimony of a security contractor who wasn't even present at the time of the theft.²¹ Randall Reid was arrested in Georgia for a shoplifting crime in Louisiana—a state he had never set foot in.²² Porcha Woodruff was arrested based on a 2015 photo, despite the fact that a 2021 driver's license photo was available.²³ And so on.

Policing errors leading to the arrest of the wrong person happen every day, and will probably continue whether or not facial recognition is used.

Besides, police have made hundreds of thousands of facial recognition searches, so the error rate of the technology is minuscule.²⁴ In fact, the error rate dropped to 0.08 percent—a fifty-fold decrease between 2014 and 2020—according to studies by the National Institute of Standards and Technology.²⁵

Facial recognition AI, if used correctly, tends to be accurate because there is little uncertainty or ambiguity in the task. Such AI is trained using vast databases of photos and labels that tell it

whether or not any two photos represent the same person. So, given enough data and computational resources, it will learn the patterns that distinguish one face from another. Facial recognition is different from other facial analysis tasks such as gender identification or emotion recognition, which are far more error prone. ^{26,27} The crucial difference is that the information required to identify faces is present in the images themselves. Those other tasks involve guessing something about a person—their gender identity or emotional state—based on their face, which puts an inherent limit on their accuracy.

Civil rights advocates have often lumped together facial recognition with other error-prone technologies used in the criminal justice system, like those that predict the risk of crime—despite the fact that the two technologies have nothing in common and the fact that error rates differ by many orders of magnitude. (The majority of people who are labeled "high risk" by predictive AI do not in fact go on to commit another crime.)

The biggest danger of facial recognition arises from the fact that *it works really well*, so it can cause great harm in the hands of the wrong people. Kashmir Hill, in her book *Your Face Belongs to Us*, details many harmful ways in which it has been used. For example, oppressive governments can and do use it to identify people in peaceful protests and retaliate against them. ²⁹

Facial recognition can also be abused by private companies. Madison Square Garden is a famous venue for sports events and concerts in New York City. In 2022, lawyer Nicolette Landi was denied entry to a Mariah Carey concert at the venue. Her boyfriend had bought the nearly \$400 tickets for her birthday. She was one of many lawyers turned away from various events at Madison Square Garden. The reason? The company that

operates the venue had banned all lawyers who worked at firms that had sued it—even if they weren't responsible for the lawsuit, and even if they were longtime visitors with season tickets. The ban was enforced using facial recognition.

When critics oppose facial recognition on the basis that it doesn't work, they may simply try to shut it down or shame researchers who work on it. This approach misses out on the benefits that facial recognition has brought. For example, the Department of Homeland Security used it in a three-week operation to solve child exploitation cold cases based on photos or videos posted by abusers on social media. It reportedly led to hundreds of identifications of children and abusers. Of course, there are more mundane benefits of facial recognition as well: unlocking our smartphones or easily organizing photos into albums based on who appears in them.

To be clear, even though facial recognition can be highly accurate when used correctly, it can easily fail in practice. For example, if used on grainy surveillance footage instead of clear photos, false matches are more likely. U.S. pharmacy chain Rite Aid used a flawed facial recognition system that led to employees wrongly accusing customers of theft. False matches happened thousands of times. The company tried its best to keep the system a secret. Fortunately, law enforcement agencies were paying attention. The Federal Trade Commission banned Rite Aid from using facial recognition for surveillance purposes for five years.³²

To summarize, a nuanced approach to the double-edged nature of facial recognition would be to engage in vigorous democratic debate to identify which applications are appropriate, to resist inappropriate uses, and to develop guardrails to prevent abuse or misuse, whether by governments or private actors.

A Series of Curious Circumstances Led to This Book

In late 2019, a former researcher from an AI company reached out to Arvind out of the blue. The company is in the lucrative business of hiring automation—a business that is filled with snake oil, as we described above. The researcher explained that people at the company knew the tool wasn't very effective, in contrast to the company's marketing claims, but the company had suppressed internal efforts to investigate its accuracy.

Coincidentally, around the same time, Arvind was invited to give a public lecture at MIT. The meeting with the researcher fresh in his mind, he spoke about AI snake oil, showcasing the sketchiness of hiring automation. Encouraged by the audience's reaction, he shared his presentation slides online, thinking that a few scholars and activists might find them interesting. But the slides unexpectedly went viral. They were downloaded tens of thousands of times and his tweets about them were viewed two million times.

Once the shock wore off, it was clear to Arvind why the topic had touched a nerve. Most of us suspect that a lot of the AI around us is fake, but we don't have the vocabulary or the authority to question it.³³ After all, it's being peddled by supposed geniuses and trillion-dollar companies. But a computer science professor calling bullshit gave legitimacy to those doubts. It turned out to be the impetus that people needed to share their own skepticism.

Within two days, Arvind's inbox had forty to fifty invitations to turn the talk into an article or even a book. But he didn't think he understood the topic well enough to write a book. He didn't want to do it unless he had a book's worth of things to say, and he didn't want to simply trade on the popularity of the talk.

The second best way to understand a topic in a university is to take a course on it. The best way is to teach a course on it. So that's what Arvind did, teaming up with Princeton sociology professor Matthew Salganik. Matt had published many foundational pieces of research showing why it's hard to predict the future with AI. We'll see two of them in chapter 3. The course was called Limits to Prediction. Matt and Arvind invited the students in the course to conduct research. One of the students in the course was Sayash.

Sayash had just joined Princeton, having previously worked at Facebook. He ultimately decided to leave Facebook to obtain a PhD and pursue public-interest technology outside a tech company. He was accepted to a few computer science PhD programs. Accepted students are invited to visit the departments in person, to meet prospective collaborators and ask questions to judge whether they would be a good fit.

When visiting departments, PhD students are advised to ask questions of this sort: What is your style of advising? How much time do your students take off? What is your approach to worklife balance? These questions are important, and they can tell you how an advisor works, but not what they value and how they think. A far more revealing question is "What would you do if a tech company files a lawsuit against you?" The answer can tell you the advisor's stance on Big Tech, how they view the impact of their research, and what they would do in a crunch. It is also unusual enough that potential advisors wouldn't have prepared their answers in advance.

Sayash asked every potential advisor this question. It carried the element of surprise, yet the scenario it described was not completely unthinkable. When Arvind answered, "I would be glad if a company threatened to sue me for my research, because

that means my work is having an impact," Sayash knew he had found the right program.

In the course on limits to prediction, students in the class were interested in predictive AI: in any and all attempts to predict the future using data, especially in social settings, ranging from civilizations to social media. Some interesting questions we looked at were: Can we predict geopolitical events such as election outcomes, recessions, or social movements? Can we predict which videos will go viral?

What we found was a graveyard of ambitious attempts to predict the future. The same fundamental roadblocks seemed to come up over and over, but since researchers in different disciplines rarely talk to each other, many scientific fields had independently rediscovered these limits. We were alarmed by the contrast between the weight of the evidence and the widespread perception that machine learning is a good tool for predicting the future.

The course included many case studies, including Google Flu Trends. This was a project that Google launched in 2008 to predict flu outbreaks by analyzing the search queries that its millions of users make every day. An increase in searches for flu-related terms could be indicative of an imminent outbreak. Google heavily promoted it as an example of AI and mass data collection used for social good. But within a few years, the accuracy of the predictions dropped precipitously. One reason was that it is hard to distinguish between media-driven panic searches and actual increases in flu activity. Another was that Google's own changes to its app changed people's search patterns in ways that weren't accounted for by the AI. Google Flu Trends ultimately ended up as a cautionary tale.³⁴ The lesson is that even in cases where it is possible to make somewhat accurate forecasts, it is very easy to get the details wrong.

Sayash found that the course confirmed his previous experiences at Facebook, where he saw how easy it was to make errors when building AI and to be overoptimistic about its efficacy. Errors could arise due to many subtle reasons and often weren't caught in testing, but only when AI was actually deployed to real users. Sayash decided to choose the limits of AI as his research topic.

After four years of research, separately and together, we're ready to share what we've learned. But this book isn't just about sharing knowledge. AI is being used to make impactful decisions about us every day, so broken AI can and does wreck lives and careers. Of course, not all AI is snake oil—far from it—so the ability to distinguish genuine progress from hype is critical for all of us. Perhaps our book can help.

The Al Hype Vortex

Since we started working together, we've come to better appreciate why there is so much misinformation, misunderstanding, and mythology about AI. In short, we realized that the problem is so persistent because researchers, companies, and the media all contribute to it.

Let's start with an example from the research world. A 2023 paper claimed that machine learning could predict hit songs with 97 percent accuracy.³⁶ Music producers are always looking out for the next hit, so this finding would have been music to their ears. News outlets, including *Scientific American* and Axios, published pieces about how this "frightening accuracy" could revolutionize the music industry.^{37,38} Earlier studies had found that it is hard to predict if a song will be successful in advance, so this paper seemed to describe a dramatic achievement.

Unfortunately for music producers, we found that the study's results were bogus.

The method presented in the paper exhibits one of the most common pitfalls in machine learning: data leakage. This means roughly that the tool is evaluated on the same, or similar, data that it is trained on, which leads to exaggerated estimates of accuracy. This is like teaching to the test—or worse, giving away the answers before an exam. We redid the analysis after fixing the error and found that machine learning performed no better than random guessing.

This is not an isolated example. Textbook errors in machine learning papers are shockingly common, especially when machine learning is used as an off-the-shelf tool by researchers not trained in computer science. For example, medical researchers may use it to predict diseases, social scientists to predict people's life outcomes, and political scientists to predict civil wars.

Systematic reviews of published research in many areas have found that the *majority* of machine-learning-based research that was re-examined turned out to be flawed.³⁹ The reason is not always nefarious; machine learning is inherently tricky, and it is extremely easy for researchers to fool themselves. Overall, research teams in more than a dozen fields have compiled evidence of widespread flaws in their own arenas, unaware that they were all part of a far-reaching credibility crisis in machine learning.

The more buzzy the research topic, the worse the quality seems to be. There are *thousands* of studies claiming to detect COVID-19 from chest x-rays and other imaging data. One systematic review looked at over four hundred papers, and concluded that *none* of them were of any clinical use because of flawed methods.⁴⁰ In over a dozen cases, the researchers used a training dataset where all the images of people with COVID-19 were from adults, and all the images of people without COVID-19

were from children. As a result, the AI they developed had merely learned to distinguish between adults and children, but the researchers mistakenly concluded that they had developed a COVID-19 detector.

We ourselves discovered flaws in many studies, mainly in the field of trying to predict civil wars (in short: it doesn't work). When we tried to publish a paper about an entire body of research being flawed, no journal was interested. It is notoriously hard to correct flaws in the scientific record. We eventually published our paper, but only after reframing it to be more palatable, as a guide to future researchers to avoid these pitfalls.

These days, when we find flawed machine learning papers, we don't even try to correct the record. The system doesn't work. In fact, in many fields, studies that fail attempts at replication by other research groups are cited *more* than those that replicate successfully.⁴¹ The party line among scientists is that science "self-corrects," meaning that the normal process of science is sufficient to root out flawed research, but everything we've seen about the process suggests otherwise.

To be clear, incorrect machine learning claims in research papers usually don't result in broken AI products on the market. If a music producer tried to predict hits using a flawed method, they would quickly find out that it doesn't work. (Commercial AI snake oil usually results from companies knowingly selling AI that doesn't work, rather than they themselves being fooled.) Still, the ocean of scientific misinformation damages the public understanding of AI, because the media tends to trumpet every purported breakthrough.

There are rays of hope, though. In summer 2022, we organized a day-long online workshop to discuss the spate of flawed machine-learning-based science. To our surprise, hundreds of scientists showed up. Based on the workshop, we led a team of

about twenty researchers across many disciplines to devise an intervention: a simple checklist that helps scientists better document how they use machine learning, which can help minimize errors and make them easier to spot when they do creep in.⁴² It's still early days, and it remains to be seen if our intervention will be adopted. At any rate, scientific practice changes glacially, and it is likely that things will continue to get worse for a while before they get better.

Let's turn to companies. While overhyped research misleads the public, overhyped products lead to direct harm. To study this, we teamed up with colleagues Angelina Wang and Solon Barocas and investigated uses of predictive AI in industry and government.⁴³ We documented about fifty applications spanning criminal justice, healthcare, welfare allocation, finance, education, worker management, and marketing. Most of these deployments are recent. During the 2010s, predictive AI extended its tentacles into many spheres of life, judging us and determining our opportunities in life based on data covertly collected about us.

We realized that while vendors of these tools aggressively chase clients, they are rarely transparent about how well their products work, or if they work at all. Notably, as far as we know, no hiring automation company has ever published a peer-reviewed paper validating its predictive AI, or even allowed an external researcher to evaluate it. Two of the leading companies made a show of external audits: Pymetrics contracted with a leading research group from Northeastern University, and HireVue contracted a noted independent auditor. But in both cases, the researchers were allowed to analyze only whether the AI was biased with respect to race or gender, and not whether it worked. The companies cleverly used a concern about discrimination to their advantage. If what you have is an

elaborate random number generator that works equally poorly for everyone, it's easy to make it unbiased!

Here, too, there are slivers of good news. Regulators are wising up to the fact that many predictive AI products don't work. In 2023, the U.S. Federal Trade Commission (FTC) warned companies that "we're not yet living in the realm of science fiction, where computers can generally make trustworthy predictions of human behavior. Your performance claims would be deceptive if they lack scientific support or if they apply only to certain types of users or under certain conditions." The key word here is "deceptive"; the FTC is authorized by Congress to police deceptive practices by companies. We hope companies will heed that warning.

If researchers and companies kindle the sparks of hype, the media fans the flames. Every day we are bombarded with stories about purported AI breakthroughs. Many articles are just reworded press releases laundered as news.

Of course, with the media so reliant on clicks and newsrooms so cash strapped, this is no surprise. Still, there are systemic problems in the industry besides crumbling revenue. Many AI reporters practice what's called access journalism. They rely on maintaining good relationships with AI companies so that they can get access to interview subjects and advance product releases. That means not asking too many skeptical questions.

Claims of AI being conscious have proven particularly irresistible to the media. When a Google engineer claimed in June 2022 that the company's internal chatbot had become sentient (and faced "bigotry"), just about every publication ran with that headline. ⁴⁵ The same thing happened when Bing's chatbot claimed to be sentient in early 2023. That's despite the fact that most AI researchers don't think there is any scientific basis for these claims.

There are many AI journalists who rise above the fray and do excellent investigative work. But so far they are a handful, constantly swimming against the tide. We've had the opportunity to discuss the hype problem with journalists and speak at a few journalism conferences. We learned about many ongoing efforts to improve the quality of tech journalism.

For example, the Pulitzer Center funds a network of journalists to work on "in-depth AI accountability stories that examine governments' and corporations' uses of predictive and surveillance technologies to guide decisions in policing, medicine, social welfare, the criminal justice system, hiring, and more." Many notable investigations have resulted from this program, including one by Ari Sen and Derêka K. Bennett for the *Dallas Morning News*. Sen and Bennett looked into Social Sentinel, an AI product used by schools across the United States to scan students' social media posts, purportedly to identify safety threats, but often misused to surveil student protests. 47

The Pulitzer Center fellowships support only ten journalists per year. In the long run, whether or not journalism can serve as a check against Big Tech's power will depend on whether funding models like these—that don't rely on clicks—can be scaled up.

AI experts have a responsibility to speak up against hype, whether it comes from researchers, companies, or the media. We are trying to do our part. In our newsletter, at AISnakeOil.com, we comment on new developments in AI and help readers separate the milk from the froth.⁴⁸

What Is AI Snake Oil?

In the late nineteenth and early twentieth centuries, snake oil peddlers were rampant in America, selling miracle cures and health tonics under false pretenses. Figure 1.1 shows a typical

INTRODUCTION 2'

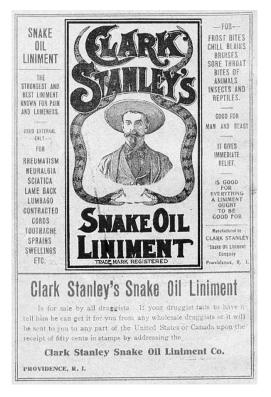


FIGURE 1.1. A 1905 advertisement for snake oil.

(Sources: https://www.nlm.nih.gov/exhibition/ephemera/medshow.html, attributed to Clark Stanley's Snake Oil Liniment, True Life in the Far West, 200 page pamphlet, illus., Worcester, Massachusetts, c. 1905, 23 × 14.8 cm. https://commons.wikimedia.org/w/index.php?curid=47338529.)

advertisement. Snake oil sellers exploited people's unscientific belief that oil from snakes had various health benefits, and their inability to tell effective treatments from useless ones. Besides, most of the concoctions being sold as snake oil didn't in fact contain any. In some cases, these medicines were ineffective but harmless. In others, they led to the loss of life or health. Until

the Food and Drug Administration (FDA) was established in 1906, there was no good way to keep snake oil salesmen accountable to their promises regarding the contents, the efficacy, or the safety of their products.

AI snake oil is AI that does not and cannot work, like the hiring video analysis software that originally motivated the research that led to this book. The goal of this book is to identify AI snake oil—and to distinguish it from AI that can work well if used in the right ways. While some cases of snake oil are clear cut, the boundaries are a bit fuzzy. In many cases, AI works to some extent but is accompanied by exaggerated claims by the companies selling it. That hype leads to overreliance, such as using AI as a replacement for human expertise instead of as a way to augment it.

Just as important: even when AI works well, it can be harmful, as we saw in the example of facial recognition technology being abused for mass surveillance. To identify what the harm is and how to remedy it, it is vital to understand whether the problem has arisen due to AI failing to work, or being overhyped, or in fact working exactly as intended. Harm and truthfulness are the two axes in figure 1.2. In this book, we're interested in everything except the bottom left part of the figure, which is AI that both works and is benign.

With this picture in mind, here's a roadmap of the rest of the book.

Chapter 2 is about automated decision-making, which is one area where AI, specifically predictive AI, is increasingly used: predicting who will commit a crime, who will drop out of school, and so forth. We'll look at many examples of systems that have failed and caused great harm. In our research, we've identified a recurring set of reasons these failures keep happening—reasons that are intrinsic to the use of predictive

INTRODUCTION 29

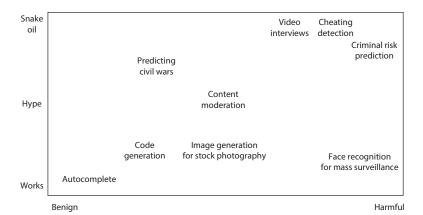


FIGURE 1.2. The landscape of AI snake oil, hype, and harms, showing a few illustrative applications.

logic in these high-impact systems. We'll end the chapter by asking if it is possible to reimagine decision-making without predictive AI, and we'll discuss what sorts of organizational and cultural adaptations we'll need in order to embrace the unpredictability inherent to consequential decisions.

In chapter 3 we'll take a step back to understand why predicting the future is so hard. Our answer is that its challenges are ultimately not about AI, but rather the nature of social processes; it is inherently hard to predict human behavior, and we'll see many reasons for this. We'll review evidence from many efforts to predict the future, from crime to children's life outcomes. We'll draw from academic studies as well as the rare cases where commercial products have been subjected to independent scrutiny. We'll look at prediction of both positive outcomes, such as succeeding at a job or publishing a bestseller, and negative outcomes, such as failing to pay back a loan; all of these turn out to be hard to predict. We'll also look at less

consequential but more easily analyzed prediction tasks such as identifying which social media posts will go viral. And finally, in addition to outcomes about individuals, we'll look at macrolevel predictions such as the evolution of pandemics. Across all of these domains, strikingly common patterns emerge, which lead us to conclude that the limitations of predictive AI won't go away in the foreseeable future.

It's simple to state the primary limitation of predictive AI: it's hard to predict the future. But with generative AI, to which we turn next, things are more complicated. The technology is remarkably capable, yet it struggles with many things a toddler can do.⁴⁹ It is also improving quickly. So, to understand what the limitations are and have some sense of where things might be going, it's important to understand the technology. In chapter 4, we hope to demystify how generative AI works.

We'll also discuss the many harms that arise from generative AI. In some cases, harms arise because the product is flawed. For instance, software that claims to detect AI-generated essays doesn't work, which can lead to false accusations of AI-based cheating. In other cases, harms arise because the product works well. Image generators are putting stock photographers out of jobs even as AI companies use their work without compensation to build the technology. Of course, there are many applications of generative AI that both work well and are broadly beneficial, such as automating some parts of computer programming (although, even here, there are minor risks that programmers should watch out for, such as the possibility of bugs in AI-generated code that might give hackers an advantage). Given the focus of the book, we won't spend much time on these beneficial applications. But we should emphasize that we are excited about them and about the potential of generative AI in general.

INDEX

Page numbers in *italics* refer to figures and tables.

access journalism, 25 account suspension, unaccountable, accuracy: news media, misleading the public, 248-50; predictions, accuracy of, 67-69; transparency, 240-41 advanced AI, existential threat: accelerating progress, 165–68; AI safety meme, 170; concluding thoughts, 177-78; expert's opinions, 151-56; global ban on, 172-74; introduction, 150-51; ladder of generality, 156-62, 160, 162–65, 163; notable historical computers, 157-59; reframe the issue from intelligence to power, 171; responses to, 172-74; rogue AI, 168-71; specific threats, defending against, 174-77. See also Artificial General Intelligence (AGI) Affordable Care Act (2010), 53 Afghanistan, hate speech (Facebook), Age of AI, The (Huttenlocher, Kissinger,

Schmidt), 251-53

moderation, 181; content moderation, seven shortcomings for, 216-18; directions for change, 33-34; evading content moderation, 198-201; for image classification, 75, 127-29; journalism, 251; safety community, 156, 170, 172, 173; safety meme, 170; seven reasons why content moderation is hard, 217; for surveillance, 127-29; for translation, 191 AI, Algorithmic, and Automation Incidents and Controversies Repository, 38 AI, future of: AI snake oil, appealing to broken institutions, 261–65; embracing randomness, 265-68; generative AI, 258-60; Kai's world, 281-84; Maya's world, 285-89; overview of, 33-34; predictive AI, 261 AI and the future of work: automation.

277-78; automation paradox, 277;

cloud computing, 276; copywriters

agency, 69-70, 247, 253, 279

AI: agents, 163-64, 169; for content

344 INDEX

AI and the future of work (continued) 261-65; definition of, 2, 26-34; and translators, 277; generative AI, embracing randomness, 265-68; 276-77; Hollywood actors and hype and harms, 29 writers, 279; National Eating Disor-AI Snake Oil: advanced AI, existential ders Association, 277; robot tax, threat of, 150-78; conclusion. 280; scriptwriting, 279; unions and 289-90; the future, predicting, workers collectives, 279; Universal 60-98; generative AI, 99-149; Basic Income, 279-80 going forward, 258-90; introduc-AI community, culture and history of tion, 1-35; myths about AI, 227-57; hype: AI research, corporate fundpredictive AI, 36-59; social media, ing of, 236-37; perceptron, 235; 179-226 scientific understanding, lack of, AI-generated: QR code image, 100; text, 237-39; springs (peaks), 235; 262-63; videos, 142-43; voices, 142 winters (valleys), 235 AISnakeOil.com, 26, 290 AI hype: AI community, culture and Alexander, Khalid, 265 history of hype, 235-39; AI hypes AlexNet, 113 and harms, illustrative applications, Algorithmic, and Automation 29; cognitive biases, 255-57; Gart-Incidents and Controversies ner hype cycle, 232; hype cycle, Repository, 38 229-30; hype vortex, 21-26; media, algorithms: advanced AI, 161-62; 25-26; previous technology hype, crime prediction algorithm, 249-50; different from, 231-35; public definition of, 39; generative AI, figures spreading, 251-55. See also AI 109-10, 114, 116, 117, 133; predicting myths, why do they persist the future, 60-61; predictive AI, AI Incident Database, 38 38-40, 48-49; social media, 207, AI myths, why do they persist: AI 219, 220-21 community, culture and history of algospeak, 201 alignment, 172-73, 174-75 hype, 235-39; AI hype, 231-35; Bing chat, misleading news headlines, Allegheny County (Pennsylvania), 249; cognitive biases, 255-57; Gart-52-53 ner hype cycle, 232; introduction, Allegheny Family Screening Tool, 227-31; news media, 247-51; public 52-53, 55 figures, 251-55; reproducibility Allstate, use of predictive AI, 10-11 crisis, 241-47; robots featured in Altman, Sam, 274 news media, 248; transparency, lack Amazon, AI-generated books, 6 of incentives for, 239-41 Amazon Mechanical Turk, 111, 115 AI research: corporate funding, American Edge, 275 236-37; recurring problems, 238; Amnesty International, 189 reproducibility crisis, 241-47 anchoring bias, 256-57 AI snake oil: 1905 advertisement for, Anderson, Wes, 7 27; appealing to broken institutions, Animal Farm (Orwell), 82

For general queries, contact info@press.princeton.edu.

INDEX 345

annotators, 115, 144-45 Anthropic, 137, 151, 260, 269 Apple, 116, 286 Artificial General Intelligence (AGI): AI safety community, 156; cognitive bias, 153; definition of, 150-51; edge cases, 154; forecasting tournament, 154-55; ladder of generality, 160; risk, estimating probability of, 155; selection bias, 152-53 Artificial Intelligence Act (AIA), 271 artificial intelligence (AI), introduction to: AI hype vortex, 21-26; ChatGPT, 3-6; dawn of, 3-6; definition of, 26-34; definition of, humorous, 14; entertainment, AI shakes up, 7-9; examples of, 14; Axios, 21 facial recognition, 15-17; generative AI, 3-6; hype and harms, illustrative applications, 29; introduction, 1-3; labeling, 12-17; labor displacement, 14; predictive AI, 2, 3, 9-12; schools and colleges, 8-9; series of curious circumstances, 18-21; snake oil advertisement (1905), 27; targeted readers (AI Snake Oil), 34-35 94-95,96 artificial superintelligence, 151 ArtStation, 126, 126 arXiv, 162-63 Asimov, Isaac, Foundation, 90, 91 Associated Press, 210, 264 AT&T, 259 ATMs, 277 attorneys, AI use, 102 autocomplete, 29, 132 autogenerated language, 140 automated decision-making, 28-29 automated hiring tools, 18, 24, 46-47 Biden, Joe, 150 automatic translation, 191 automating bullshit: autogenerated language, 140-41; bullshit, For general queries, contact info@press.princeton.edu.

definition of, 139; ChatGPT, 139, 140; CNET, 140; defamation by chatbot, 140; examples of, 139-41 automating script writing, 8 automation: AI and the future of work. 277-78; automation bias, 255; content moderation and, 181; hate speech, directed at Black people, 186; hiring automation, 18, 24, 66; labor concerns, 276-77, 278; low-wage workers, effect on, 280; positive effects, 277-78; robot tax, 280 automation bias, 50-51, 255 automation paradox, 277 automation's last mile, 278 autonomous driving, 13, 14, 153-54 backpropagation, 109, 113

Balenciaga Pope, 7 Bankman-Fried, Sam, 234 Bard (renamed Gemini), 5, 135, 136-37 Barocas, Solon, 24, 42-43 base model, 133-34 basic reproduction number (COVID), Be My Eyes (app), 99-100 benchmark datasets, 112, 237-38, 241 benchmarking, 112, 116 Bender, Emily, 248 Bennett, Derêka K., 26 Better Business Bureau, 213 biases: AI lack of, 79; benchmarks, 116; COMPAS, 79, 80; criminal risk prediction tools, 11; image generators, 103; of policing, 80; selection bias, 152. See also cognitive biases Bing's chatbot, 25, 101, 249 biological risk, 176-77

Bitcoin, 233-34

346 INDEX

Bitcoin mining, 234 Black, Rebecca, Friday (song), 88 blackout challenge, 219 blasphemy, 192-93 Bledsoe, Drew, 82 Bloomberg, 249 Bosnia, content moderation failures, 190 bots: AI agents, 163-64; benchmarks, 241; companion bots, 102; generative AI, 101; rock-paper-scissors, 137; scraping, 115. See also chatbots Brady, Tom, 82 broken institutions: AI-generated text, 262; definition of, 263; educational institutions, 262-63; efficiency as a selling point, 263; gun violence and, 263; journalism, 262; law enforcement, 263-64; predictive AI, 265; snake oil, demand for, 261-62; state of hiring, 261-62; weapons, using AI for detecting, 263 Broward County (COMPAS), 79-80 brute-force intervention, 74, 131 bug finding tools, 175-76 bullshit: as autogenerated language, 140-41; ChatGPT, 139, 140; CNET, 140; defamation by chatbot, 140; definition of, 139; examples of, 139-41; generators of, 197 Butterfly Effect, 63

Canada, 128–29, 280
cancer, 69, 239
Center for AI Safety, 152
Centers for Disease Control and
Prevention (CDC), 92, 93, 215
Chai, 102–3
chance events, 75
channels, automated notifications of
copyright claims, 209

Charter, 259 chatbots: arXiv, 163; factual information weakness, 5; inappropriate uses, 6; introduction, 2; National Eating Disorders Association, 277 chatbots, generative AI: Bard (renamed Gemini), 137; base model, 133-34; bullshit, automating, 139-41; Claude, 137; defamation by chatbot, 140; fine tuning, 134-35, 137-38; inaccurate outputs, 147-48; inappropriate outputs, 101-2; internal representations, 138-39; limitation of, 136-37; a meta task, fine-tuned for, 135-36; Othello games, 138; philosophical dimension, 137; practical dimension, 137; translation tool, 134-35 ChatGPT: autocomplete, 132; automating bullshit, 139, 140; Be My Eyes, 100; computer programmers, 4-5; educator's curricula, 262-63; fine tuning, 135-36; G stands for generative, 133; humor, detecting, 185; introduction, 3-6; ladder of generality, at time of writing, 163; misinformation, 102; P stands for pretrained, 135; schools and colleges, 9; T stands for Transformer, 131; text generation, 131; token, generation of, 133 Chattanooga, Tennessee, 260 checkers-playing program, 107

Charlie Bit My Finger, 86-87

chess-playing computer, 108 Chiang, Ted, 278 Chicago (ShotSpotter), 264 Chicago Tribune, 1948 election, 56 child sexual abuse, 194–95 child sexual abuse imagery, 194

INDEX 347

Children's Online Privacy Protection commercial AI snake oil, 23 Act (COPPA), 282 Common Crawl, 114 community networks, 259-60 China: AI, oversight of, 271; chatbots, use of, 271; Core Socialist Values, Community Notes, 220 271; COVID deaths, unpredictabilcompanion chatbots, 102-3 ity of, 95-97; Cybersecurity Law, COMPAS, 41, 55, 79-80, 239 271; facial recognition, 127; misincomputational predictions, effectiveformation, 193; social media apps ness of, 67 banned, 216 computer vision, 111-13 civil wars, 23, 189, 243-44, 250-51 computers, notable historical, 157-59 classifier: child sex abuse, 194-95; computers, predicting the future with, COVID vaccines, 196; image recognition classifier, 173; leakage Consumer Finance Protection Bureau, example, Russian and American 270 tanks, 244; movie review sentiment Content ID: abuse by police officers, classifier, 134; slur words, discerning, 209; algorithm, 207; fair use con-186; suicide or self-injury, immicept, 208; overview of, 206-7; nent, 203, 204; text generation, 129 running amok, 207-8 Claude (chatbot), rock-paper-scissors, content moderation: amplification of problematic content (Zuckerberg), 221; classifiers, 184-85; a dead end, Clean Air Act, 269 223; evading, 198-201; future of, Clean Water Act, 269 Clearview AI, 127, 128–29 223-26; hate speech, 186, 222; misclickbait, 248, 249 takes, 183-84; problem of their own cliodynamics, 91-92 making (companies), 218-23; procloud computing, 276 cess of, 181-83, 182; seven shortcom-CNET, 140, 262 ings of AI, 216-18, 217; social media's problem, 218-23; stages of AI-based CNN, 247 cobra population, reducing, 46 content moderation cycle, 214; Code Red, 176 subreddits, 224-25, 226 Codex, 245 content moderators, 182, 182, 185, cognitive bias, 153, 255-57 190-91, 214-15, 223 cognitive biases: anchoring bias, context, discerning, 183-88 256-57; conclusion, 257; explanatory Cook County (Illinois), 52 depth, illusion of, 255-56; Future of Correctional Offender Management Life Institute, 256; halo effect, 256; **Profiling for Alternative Sanctions** humans, 231; illusory truth effect, (COMPAS), 41, 55, 79-80, 239 counter notice, 208 256; overview of, 231, 255; priming, 256; quantification bias, 257 COVID interventions, case study

(New Zealand), 96

Comcast, 259

348 INDEX

COVID pandemic prediction: basic Damon, Matt, 234 reproduction number, 94-95; DARPA, 259 COVID interventions, case study data annotation, 111, 145, 146 (New Zealand), 96; deaths, data centers, surveilling, 173 unpredictability of, 95-97; flu, data collection: Google, 20; governdifference between, 94; overview ments, 77; NSA, 78; predictive AI of, 93; short-term forecasting, companies, 45; tech companies, 77 93-94; strategic decisions, 98 data leakage, 22 COVID-19: detector, 22–23; lab-leak datasets: AI research, 167; benchmark theory, 177; machine learning, datasets, 112, 237, 241; ImageNet, 114; ImageNet competition, 121-22; 22-23; misinformation, 195-96 creative labor, generative AI predictive AI, 41-42; scraped appropriates, 122-26 criminal justice, 40-41, 78-81, 266 criminal risk assessment, 68-69 criminal risk prediction, 11, 51-53, 59, David, Larry, 234 66, 256 criminal risk prediction systems: Allegheny Family Screening Tool, 52-53, 55; Ohio Risk Assessment System, 51; Public Safety Assessment, 51, 52 criminal risk prediction tools, 11 criti-hype, 178, 253, 254 cryptocurrency, 233-35 cultural: competence, 192; incompetence, 188-94; products, 83-84, 98 culture: AI Images, 126; equitable distribution of consumption, 83; and history of hype, 235-39; scams, 142 ImageNet, 114-18; information defense in depth, 176 security, 176; memes, 87; Reddit, 225 Deng, Jia, 111 cumulative advantage, 83-84, 85. See also rich-get-richer dynamics Custodians of the Internet (Gillespie), 210 270 cybersecurity, 175-76, 177

Dallas Morning News, 26 Dall-E, 2, 7, 121, 123 D'Amelio, Charli, 88

datasets, downside to, 115-16, 122; social datasets, 76; sociological datasets, 71; Stable Diffusion, 122 deceptive practices, 25, 273 deep learning: algorithms, 117; computer vision, 113; image classification, 119; ImageNet Challenge, 74, 112-13; ladder of generality, 161-62, 162, 163; radiologists, replacing, 238-39; text generation, 129. See also ImageNet, technical and cultural significance of deep neural networks, 109, 114, 118, 120 deep synthesis systems, 271 deepfakes, fraud, malicious uses: AIgenerated videos, 143; AI-generated voices, 142; deepfakes, 142-43, 148; Department of Commerce, 270 Department of Homeland Security, 17, Diary of a Young Girl, The (Frank), 82 diffusion model, 121, 122 digit recognition, 109, 110 digit recognition machines, 109 digital infrastructure, 259, 260

INDEX 349

Digital Markets Act (DMA), 270–71 Digital Services Act of 2023 (DSA), 223–24, 270 disease prediction, 92 Disney, 83 dog whistles, 191 DoorDash, 213 Douek, Evelyn, 224 downranking or demotion, 183

EAB Navigate, 37, 38 earthquakes, 68 edge cases, 154 EdTech, 286 eight billion problem, 76, 97 election forecasting, 56 ELIZA, 165-66 ELIZA effect, 166 emotion recognition, 16 entertainment, 7-9, 88, 89, 148, 226 Epic, 227, 239 Epic sepsis model, 227–29, 230 Ethiopia, 189, 191, 260 Etsy, 213 EU, 260, 270-71 Evans, Benedict, 276 Executive Office of the President, 270 existential risks, 31, 151, 154, 255. See also advanced AI, existential threat expert systems, 161, 236

Facebook: American Edge, 275; Black people, getting locked out, 187; community standards document, 181, 195, 196; content moderation failures, 188; content moderation rules, 185–86; content moderators, 190–91, 191–92; hate speech, 189, 190, 218; Koobface, 199; Middle East, takedowns of terrorist content, 190;

oversight board, 212; Phan Thi Kim Phuc (Napalm Girl), 210, 211; post removal, 182, 182-83; regulatory capture, 274; Rohingya (Myanmar), persecution of, 189; role in society, 179; Sri Lanka, hate speech, 190; suicide prevention, 202-3; suppression of conservative political movements, 187; Tigray civil war, role in, 189; welfare checks, 203, 204; worst of the worst investigation, 187 facial analysis technology, physical billboards, 128 facial recognition: accuracy of, 15-16; banning, 129; benefits of, 17; China, 127; Clearview AI, 127, 128; danger of, 16; false arrests, 15; introduction, 15-17; police officers abuse of, 127-28; racial bias, 15; Rite Aid, falsely accusing customers, 17; surveillance, abusive use of, 127-29 fair use provision (U.S. copyright law), 123, 208 false dichotomy, cutting through, 268-74 Faulkner, Judith, 228 Federal Trade Commission (FTC), 17, 25, 140 feedback loops, 97-98, 108-9, 268 Ferguson, Viana, 186 FICO (Fair, Issac, and Company), 66 FICO score, 66 fine-tuning, 118, 134, 135-36, 137-38, fingerprint matching, 194, 206, 218 Finland (UBI), 280 Firefly (Adobe), 7 First Amendment (United States Constitution), 193, 272 flawed democracy (India), 216

350 INDEX

Flu Trends, 20 FluSight competition, 92, 93 Food and Drug Administration (FDA), 28, 270 Forrester, Jay, 64 Foundation (Asimov), 90 4chan, 142, 180 Fragile Families Challenge, 70-73, 97 Fragile Families Challenge, disappointment in the end, 73-78 Frank, Anne, The Diary of a Young Girl, Frankfurt, Harry, 139 Friday (Black), 88 FTX, 234 the future, getting specific: cancer, 69; criminal risk assessment, 68-69; earthquakes, 68; genetic diseases, 69; life outcomes, 69-70; people's futures, 67; phenomena can predict, 67; phenomena can't predict, 67;

predictions, judging accuracy of, 67–69; the weather, 67–68 Future of Life Institute, 151, 254, 256

gaming, 47–48
gaming, AI incentivizes, 46–48
Gangnam Style (PSY), 88
Gartner hype cycle, 231–33, 232
gatekeepers, 88–89
Gebru, Timnit, 102
Gemini, 5, 135, 136–37
gender identification, 16
General Data Protection Regulation
(GDPR), 270
Generation Alpha, 282
generation time, 95
generative AI: AI-generated image,
QR code, 100; appropriating

creative labor, 122-26, 124-25, 126;

automating bullshit, 139-41; cost of improvement, 143-46; deepfakes, fraud, malicious uses, 142-43; failure and revival, 107-10, 110; functional QR code, 100; historical background, 105-7, 106; image classifying and generating, 118-22, 119, 120, 121, 122; from images to text, 129-33; introduction, 2; in large-scale labor, 143-46; models, 123, 143, 260; from models to chatbots, 133-39; overview of, 3-6, 30, 99-105; for surveillance, 127-29, 147; taking stock, 146-49; text-to-image generation, 7-8; training machines to see, 111-13 genetic diseases, predicting, 69 Getty Images, 123 Ghost Work (Gray, Suri), 278 Gig City, 260 Gillespie, Tarleton, Custodians of the Internet, 210-11 GitHub Copilot, 5 Google: account suspension, unaccountable, 213-14; chatbot (Bard), 5, 137; child sexual abuse, classifier for, 194-95; dataset sizes, 114; internal criticism, silencing, 102; Mark and his toddler, 183-84, 195, 212; public knowledge to trade secrets, shift from, 260; Viacom lawsuit, copyright violation, 206; voice assistant, 248 Google Cloud, content moderation mistake, 183-84 Google Drive, 276 Google Flu Trends, 20 Google Photos, 115-16 Google Translate, 191 GPT-3, 144

INDEX 35

GPT-4: alignment, negating the effect of, 175; human-level performance, 241; hype by public figures, 254; limitation of, 136; professional exams, 276 gradient descent, 117, 162 graphics processors (GPUs), 113, 131 Gray, Mary L., Ghost Work, 278 Grisham, John, 82 Gundersen, Odd Erik, 243

halo effect, 256 Harry Potter (Rowling), 82 hate speech: Afghanistan, 190; content moderation, 186, 222; cost of improvement, 143; Facebook, 189, 190, 218; inciting violence, 222; social media platforms, 181; Sri Lanka, 190; Tigray civil war, 189 health information, misinformation, healthcare decisions, 43-45 herding effect, 161 Herzegovina, content moderation failures, 190 Hill, Kashmir, Your Face Belongs to Us, 16 Hiltzik, Michael, 247 Hinton, Geoffrey, 109, 112-13, 238-39 HireVue, 24-25, 42, 239, 261-62 hiring automation, 18, 24, 66 Hitler, A., Mein Kampf, 142 Holocaust denial, 192 human: content moderators, 182, 185, 214-15, 223; ingenuity, 198-201; oversight, 50 human extinction. See Advanced AI, existential threat Huttenlocher, Daniel, The Age of AI, 251 hype, 235-39

Idaho (autogenerated language), 141 illusory truth effect, 256 image classification, 118-20, 127-29 image generation, 7, 12-13, 119, 121-22, 123, 124-25 image generators: biases and stereotypes, 103; Dall-E, 2, 121-22; downside of, 103; fake movie trailers, 7; introduction, 7–9; Midjourney, 2; Stable Diffusion, 2, 121-22 ImageNet, 111, 161 ImageNet, technical and cultural significance of, 114-18 ImageNet Challenge, 74 ImageNet competition, 111-13, 118, 121-22 images, generative AI: AI-generated image, QR code, 100, 100-101; Be My Eyes, 100; ChatGPT, 100; classifying and generating, 118-22; Clearview AI, 127; copied without compensation, 123, 126; Dall-E 2, 123; deepfakes, 142-43; fair use provision, 122-23; generative AI models, 123; images to text, 129-33; Lensa, 103; perceptron, 106; pornographic, 103; Stability AI, 122; Stable Diffusion, 122, 123; watermarked, 123; web scraping, 115 images, introduction to, 7, 22-23 images, myths: leakage, 244; perceptron, 235; robots, 247, 248, 256 images, social media: child sexual abuse imagery, 194-95; content moderation, 180; context, taken out of, 183-84; Facebook's internal documents, 185-86; Napalm Girl, 209-11 images to text: autocomplete, 132-33; deep learning, 129; introduction, 129-30; long-range dependencies, 130-31; processing matrices, 131

352 INDEX

incentive misalignment, 12 incitement to violence, 187, 195 India, 46, 127, 145, 190, 216 individuals to aggregates: cliodynamics, 91-92; demand forecasting, 90; disease prediction, 92; geopolitical events, predictions about, 91; limits to prediction, reasons for, 97-98; pandemic prediction (COVID), 93-97; psychohistory, 90 Indonesia, content moderation failures, 190 information security, 175, 176 inner monologue, 163 Instagram, 203, 213, 222 Institute of Advanced Study, 63 instruction following, 135-36, 144, 172-73 internet connectivity, 259 interventions: alleviating poverty, 70; COVID interventions, case study (New Zealand), 96; lotteries, use of, 268; misinformation interventions, 197; schools, 37; suicide prevention, 201-6; welfare checks, 204-5 iPad kids, 282 iPhone 13 Pro, 115 Iraq, cultural incompetence, 190 irreducible error, 69, 75, 76 Iyer, Ravi, 222, 223 James, William, 238

James, William, 238 Jelinek, Frederick, 116 job-candidate-assessment model, 239 John Carter, 83

Kai's world, growing up with AI, 281–84 Kapoor, Sayash, 19–21, 188 Karya, 145 Keller, Daphne, 224
Kenya, content moderation failures, 190
kidney transplant matching, 11–12
killer robots, 247, 256
Kim Phuc, Phan Thi, 209
Kindel, Alex, 72–73
Kissinger, Henry, *The Age of AI*, 251
Kjensmo, Sigbjørn, 243
Koobface, 199
Krizhevsky, Alex, 112–13

labels, 15-16, 119, 132, 194-95 labor displacement, 14, 277, 279 labor exploitation, 143-46, 148 labor rights, 35, 269 ladder of generality: first rungs, 160; notable historical computers, 157-59; overview of, 156-62; at time of writing, 163; until early 2010s, 162; what's next, 162-65 Landi, Nicolette, 16-17 language models, 141, 143-44 lawyers, 16-17, 212, 241, 276 leakage, 22, 244 LeCun, Yann, 153, 164 legacy admission, 287-88 Lensa, 103 Li, Fei-Fei, 111 liar's dividend, 148 life and death, a matter of, 201-5 life outcomes, predictability of, 75 life outcomes, predicting, 69-70 Lighthill, James, 235 Limits to Prediction, 19 Lipton, Zachary, Troubling Trends in Machine Learning Scholarship, 238 Live Time, 230 long-range dependencies, 130-31 long-term predictive power, 73-78

Lorenz, Edward, 63

INDEX 353

low-wage workers, 180, 276, 278, 280 Lucas, George, 82–83 luck, 81–82, 85, 86, 137, 267 Lundberg, Ian, 72–73

machine learning: algorithms, 39; child sexual abuse, 194-95; COVID-19, 22-23; criminal risk prediction, 66; definition of, 13; divorce, predicting, 70-71; the future, predicting, 60-61, 65-66; hate speech, 195; hiring automation, 66; hypertension, predicting, 44-45; intervention, 24; ladder of generality, 160; music industry, predicting hit songs, 21-22; personalized ads, serving, 66; random guessing, 22; suicide prevention, 202-3 machine-learning-based: research, 22; science workshop, 23-24 machines, 109, 111-13, 277 Madison Square Garden (Abuse of AI), 16-17 main character of Twitter, 89 Marcus, Gary, 154 Mark and his toddler, 183-84, 195, 212 Markup, The, 259 Maryland, use of predictive AI, 10–11 mass collaboration, 72 Mastodon, 179, 225-26, 260 Maya's world, growing up with AI, 285-89 McCulloch, Warren, 105 McLanahan, Sara, 71-73 media: AI hype, 25-26, 231; AI hype vortex, 21, 23, 25; Flu Trends, 20; killer robots, 256; state-sponsored influence operations, 199. See also social media media influencers, 88

Medicare, 9-10, 268 Mein Kampf (Hitler), 142 meme lottery, the, 86-90 meritocratic society, 83 Meta, 203, 269 Michigan, overautomation failure, 49 Microsoft, 5, 143, 176 Middle East, takedowns of terrorist content, 190 Midjourney, 2, 7 Minority Report, 75, 128 Minsky, Marvin, 153, 235 misinformation: AI hype vortex, 21-26; ChatGPT, 102; illusory truth effect, 256. See also automating bullshit misinformation, social media: authoritarian governments use of, 197; COVID-19, 195-96; cultural incompetence, 193; detection of, 197; health information, 197; language models, 197; machine learning, the predominant approach, 195; natural engagement pattern, 220; removal of, 197; taking down, 216 Mission Impossible: Dead Reckoning, 150, 178 Mitchell, Margaret, 102 Mitchell, Melanie, 236, 238 models: base model, 133-34; language models, 134, 196-97; translation model, 191 Mollick, Ethan, 286 Mollick, Lilach, 286 Mona Lisa replication, 124-25 Mount St. Mary's University, 36-37 music industry, 21–22 Musk, Elon, 180 Myanmar, 187-89, 260 Myth of Artificial Intelligence, The (Suchman, Whittaker), 252

354 INDEX

myths about AI: AI community, culture and history of hype, 235-39; AI hype, 231-35; AI hype, public figures spreading, 251-55; Bing chat, misleading news headlines, 249; cognitive biases, 255-57; Gartner hype cycle, 232; introduction, 227-31; news media misleads the public, 247-51; overview of, 32-33; reproducibility crisis, 241-47; robots featured in news media, 248; transparency, few incentives for, 239-41

249; Epic, 248; headlines, 249, 250; hype, underlying reasons for, 251; images of AI, 247; robots in news media, 248; universities, press releases from, 250 next-word-prediction, 132-33, 248 Nimda, 176 No to AI Images, 126, 126 noise, 121 nonproliferation, 174 notable historical computers, 157-59 NVIDIA, 114

Napalm Girl, 209-11 Narayanan, Arvind, 18-21, 62, 81, 234-35 National Center for Missing and Exploited Children (NCMEC), 194 National Eating Disorders Association (NEDA), 277 National Institute of Standards and Technology, 15 natural engagement pattern, 220, 221

Nature, 109

Netflix, 39

Netherlands, 48-49, 77, 85-86

neural networks: 5-layer neural network, 110; classify images of dogs, 119; failure and revival, 107-10; image generation, 121; ImageNet, 114; ImageNet competition, 112-13; text generation, 131. See also deep neural networks

NeurIPS, 237, 246 New York Magazine, 144 New York Times, 247, 248 news media, misleading the public: accuracy numbers, 248-51; AI hype, debunking, 247-48; civil war prediction, 250-51; clickbait, 248,

object recognition, 13 Ohio Risk Assessment System (ORAS), 51 Online Safety Bill (2022) (UK), 143 opaque AI incentivizes gaming, 46-48 opaque AI models, 46-48 OpenAI: AGI, 151; AI regulations, 274; ChatGPT, 3-5; Codex, 245; fine tuning, 135-36; public knowledge to trade secrets, shift from, 260; regulatory capture, 274; Sama, pay disparity, 144; toxic speech and offensive outputs, 269 optimization mindset, 266 Optum's Impact Pro, 54-55 Oregon experiment, 268 Ortiz, Karla, 290 Othello, 138 overautomation, 48-51 overreaching, 187, 197 Overton window of speech, 211, 212

paper clip maximizer, 2, 171 Papert, Seymour, 235 partial lotteries, 266-68, 288 the past, predicting, 194-97 peaks, (springs), 235

INDEX 355

peer review, 242 peer-reviewed, 24, 228, 239, 242 perceptron, 105-7, 106 Person of Interest, 75 Photoshop, 143 Pineau, Joelle, 246 Pitts, Walter, 105 policymaking, 212, 217, 218 political units, datasets, 91 power, defined, 171 PowerPoint presentations, 278 precarious work, definition of, 144 predicting the future: AI fails, five reasons why, 59; computers, predicting the future with, 62-67; COVID interventions, case study (New Zealand), 96; criminal justice, predictions in, 78-81; fragile families challenge, 67-70; fragile families challenge, disappointing end, 73-78; getting specific, 67-70; individuals to aggregates, from, 90-97; introduction, 60-62; limits to, 97-98; meme lottery, 86-90; predicting success, 81-86; prediction systems, 60; weather forecasting, 62-64, 65; wrong people, predictions about, 51-53 predictive AI: Allstate, 10-11; automated decision-making, 28-29; concluding thoughts, 58-59; failure, example of, 9-10; five reasons why AI fails, 59; gaming the system, 11-12, 47-48; good prediction is not a good decision, 43-45; in hiring, 10; in industry and government, 24; insurance rates, 10-11; introduction, 2, 3, 9-12, 36-38; life-altering decisions, making, 38-43; predicting the future,

29–30. See also the future, getting specific predictive AI goes wrong: automated hiring tools, 18, 24; concluding thoughts, 58-59; EAB Navigate, 37; five reasons why AI fails, 59; healthcare decisions, 43-45; incentivizing gaming, 11-12, 46-48; inequalities, exacerbating existing, 53-55; introduction, 36-38; life-altering decisions, making, 38-43; overautomation, 48-51; randomness, 56-58; university retention rate, 36-37; world without prediction, 56-58; wrong people, predictions about, 51-53 presidential elections: 1948 presidential election, 56; 2016 presidential election, 199; Chicago Tribune, 1948 election, 56 pretraining, 135, 136, 163 priming, 256 pro-ana posters, 221 products, overhyped, 24-25 programmatic interface, 175 prohibition, 275 proprietary, 228, 246-47 PSY, Gangnam Style, 88 psychohistory, fictional science of, 90, 91 Ptacek, Thomas, 3-4 public figures, spreading AI hype: Daniel Huttenlocher, 251; Eric Schmidt, 251; Future of Life Institute, 254-55; Henry Kissinger, 251 public interest, promoting AI: false dichotomy, cutting through, 268-74; introduction, 258-61; regulation, limitations of, 274-75; snake oil, appealing to broken institutions, 261-65

356 INDEX

Public Safety Assessment (PSA), 51, 52–53 Pulitzer Center fellowships, 26 Pymetrics, 24–25

QR codes, AI generated image, 100 QR codes, 100–101 qualitative criteria, 68 quantification bias, 257

racial bias, 54, 79-80 radiologists, replacing, 238-39, 276 Rahimi, Ali, 237-38 randomized controlled trials (RCTs), 45 randomness: discomfort with, 56-58; election forecasting, 56-57; embracing, 265-68; housing allocation lotteries, 58; psychohistory, 90 Recht, Benjamin, 237 recidivism-prediction models, 239 recommendation algorithms, 31, 88, 211, 220-22 recommender systems, 260, 271 Reddit, 213, 224-25, 226 Reddit-style moderation, 224-25 regulation: automated notifications of copyright claims, 209; content moderation decisions, 205-6; copyright, 206-9, 209; fake copyright strikes, 208-9; hate speech, 272; horizontal regulations, 270, 271; limitations of, 274-75; overzealous regulation, 275; platforms have free reign, 205; social media, 205-9; U.S. copyright law, 208; vertical regulations, 270, 271 regulation: cutting through the false

dichotomy: China, 271; companies,

268-69; enforcement, 273; envi-

ronmental protection, 269; EU,

270–71; food safety, 269; labor rights, 269; lagging behind, 272; myths, 269, 272–73; regulation, definition of, 268; regulation, improvement of, 273; self-regulation, 271, 272; tech regulation, 272–73; United States, 270, 272

regulations: AI for surveillance, 129; annotators, 144; China, 271; Clean Air Act, 269; Clean Water Act, 269; Digital Services Act, 223–24; enforcement of, 273; false dichotomy, cutting through, 268–74; labor rights, 269; limitations of, 274–75; U.S. National Labor Relations Board, 279

regulatory capture, 274–75
Reid, Randall, 15
relative accuracy, 79, 228
reproducibility, definition of, 242
reproducibility challenge, 246
reproducibility crisis in AI research:

AI-based science, 244–45; civil war prediction, 243–44; Codex, 245; commercial models, reliance on, 245; conclusion, 246–47; improving reproducibility, 245–46; leakage, 244; NeurIPS, 246; online workshop, 244; oversight, lack of, 231; overview of, 241–42; peer review, 242; reproducibility challenge, 246; scientific advances, genuine, 245; social psychology, 242; 2018 study, 243 resource scarcity, eliminating, 265, 268

Retorio, 47
reverse image search, 118–20, 127
rich-get-richer dynamics, 83–84, 86, 88, 266–67, 268
risk assessment tools, 40–41, 79–80
Rite Aid, falsely accusing customers, 17

INDEX 357

River Dell High School students, 290 Robodebt scandal (Australia), 49 robot lawyer, 104 robot tax, 280 robotics, 32, 256 robots, 247, 248, 256 rock-paper-scissors, 136-37 rogue AI: AI safety community, 170; AI safety meme, 170; overview of, 168-71; paper clip maximizer, 168-69; power, definition of, 171; power-seeking behavior, 169; reframe the issue from intelligence to power, 171; superintelligent beings, 171; us-versus-it argument, 168 Rohingya (Myanmar), persecution of, 188-89 Roose, Kevin, 101 Roosevelt, Franklin D., 275 Rosenblatt, Frank, 105, 153, 235 Rowling, J. K., Harry Potter, 82 rules-based system, 165-66 Russakovsky, Olga, 111 Russia, 199, 244 Sacco, Justine, 89, 222 Salganik, Matthew, 19, 62, 71-73, 76, 84-85

Sacco, Justine, 89, 222
Salganik, Matthew, 19, 62, 71–73, 76, 84–85
Sama, 144
San Diego's surveillance program, 265
Santa Clara Principles, 192
Schmidt, Eric, *The Age of AI*, 251, 252
schools and colleges, 8–9, 262, 265
Schwartz, Barry, 267
Science, 245
science self-corrects, 23
Scientific American, 21
scientific research, 242–47, 250, 267
screen time, 285
scriptwriting, 279

selection bias, 152-53 self-driving cars, 14, 32, 153, 168. See also autonomous driving Sen, Ari, 26 sentient, 25, 101 sepsis, predicting the risk of, 227-29 Seuss, Dr., 82 shadow banning, 183 Shazam, 207 ShotSpotter, 263-64 Shutterstock, 123 simplicity, 266 simulation, 64-66 Simulmatics, 65 snake oil: AI hype vortex, 21-26; AI journalism, 251; Artificial General Intelligence, 152; broken institutions, appealing to, 261-65; cognitive biases, 257; commercial AI snake oil, 23; content moderation AI, 226; criti-hype, 178; definition of, 2, 26-34; investigating accuracy, 253; Mark Zuckerberg, 179; myths about AI, 230; predictive AI, 37, 57, 58, 59, 98; robot lawyer product, 104 snake oil advertisement (1905), 27 social media: automated notifications of copyright claims, 209; content moderation, 194-97; context, taken out of, 183-88; coordinated mass harassment, 89; cultural incompetence, 188-94; drawing the line, 209-16; evading content moderation, 198–201; Facebook's moderator training materials, 196; fingerprint matching, 194, 206, 218; introduction, 179-83; life and death, a matter of, 201-5; Overton window of speech, 212; overview of, 31; the past, predicting, 194-97; problematic

358 INDEX

social media (continued) Strange World, 83 content, amplification of, 221; students: AI-generated text, 262-63; ChatGPT, 9; cheating detection problems of their own making (platforms), 218-23; regulation, software, 33-34; exercises and pedaadding to the mix, 205-9; seven gogical materials, 35; ImageNet, 111; reasons content moderation is hard, introduction, 9; Limits to Predic-217; stages of AI-based content tion, 19-20; mental health support, moderation cycle, 214. See also 264; online course, Bitcoin and content moderation cryptocurrencies, 235; partial lottersocial media companies, 32, 88, 147, ies, 267; predictive AI, 36-37; Prince-191, 216 ton, 8; River Dell High School, 290; social media, content moderation: Social Sentinel, 26, 264-65 content moderation, evading, success, predicting, 81-86 198-201; fingerprint matching, 194; Suchman, Lucy, The Myth of Artificial the future of, 223-26; machine Intelligence, 252 learning, 194-95; platform's failures, suckers list, 10 218-23; seven shortcomings of AI, suicide, a matter of life and death, 201-5 216-18 suicide prevention efforts, 201-3 social predictions, improving, 74, 75, 76 superintelligence, 151, 166, 171 social science, 70-71. See also Fragile superintelligent beings, 171 Families Challenge Support Vector Machine (SVM), Social Sentinel, 26, 264-65 109-10 South Korean government, 127 Suri, Siddharth, Ghost Work, 278 surveillance, 127-29 spam, 66, 181, 198 spam classifiers, 66 Sussman, Gerald, 153 specific threats, defending against, Sutskever, Ilya, 112–13 SVM, 109-10, 112, 114 174-77 symbolic systems, 108-9, 160, 161 speech, regulation of, 32 springs (peaks), 235 system dynamics, 64 Sri Lanka, 190, 260 systems thinking, 224 Stability AI, 7, 134. See also Stable Diffusion tech journalism, 26, 284 Stable Diffusion, 2, 7, 121–22, 122–23 teen mental health, 285 Star Wars, 82-83 Telangana (India), 127 STAT, 230 Telegram app, 142-43 Steinhardt, Jacob, Troubling Trends in Test of Time Award, 237-38 Machine Learning Scholarship, 238 Tetlock, Philip, 154 STEM majors, 37 text generation, 129-33 stimulus checks, COVID-19, 38 text-to-image generation: alreadystock photos, 104 deployed models, harms of, 254;

INDEX 359

appropriating creative labor, 122-23; Dall-E 2, 7; diffusion model, 121-22; entertainment, 7-8; Firefly, 7; Midjourney, 7; Mona Lisa, 124; Stable Diffusion, 7 text-to-video generators, 8 threats (specific), defending against, 174-77 threshold of usefulness, 166 Tigray, Ethiopia, 189, 191 TikTok, 86, 88, 218-19 tobacco companies, 239, 274 token, 132, 133 toxic content, 143-44 training data: appropriating creative labor, 122-23; asthmatic patients, 43-44; chatbots, 5; contentmoderation classifiers, 184-85; cost of improvement, 144; creating, 35; from images to text, 129, 132; leakage, 244; machine learning models, 66; from models to chatbots, 136; Ohio Risk Assessment System, 51; Optum's Impact Pro, 54; Public Safety Assessment, 52 training data, generative AI: biases and stereotypes, 103; competitions, 112; Google, 114; ImageNet, 115; Telangana, 127 translation tool, 134 transparency: Artificial Intelligence Act, 271; China, 271; content moderation, future of, 223-24; Digital Services Act, 270; language models, 141; OpenAI, 274; ShotSpotter, 264 transparency, incentives for: accuracy measurements, 240-41; benchmark datasets, 241; COMPAS, 239; early-stage startups, 239-40; Epic, 239; GPT-4 (OpenAI), 241; Hire-

Vue, 239; top-N accuracy, 240; venture capitalists, 240 Troubling Trends in Machine Learning Scholarship (Lipton, Steinhardt), 238 truth of statements, evaluating, 196, 197 Turchin, Peter, 91 Turing, Alan, 156 Turing Award, 109, 164 Turkey, 216 Twitter: account suspension, unaccountable, 213–14; bridging-based rankings, 220; content moderation, 180; content moderation mistake, 184; meme lottery, 87; Turkey, blocking opposing voices, 216 20th Century Fox, 82-83 2023 Hollywood strikes, 7–8

Uber, 213
United Artists, 82
United States: Epic sepsis model, 227;
false arrests due to facial recognition, 15; Social Sentinel, use in schools, 26
United States Postal Service, 109
United States, predictive AI: Affordable Care Act, 53; criminal risk prediction systems, 51–52; election forecasting, 56; racial disparities in policing, 55

United States, promoting public interest: AI for detecting weapons, 263; AI regulation, 270; community networks, 259–60; environmental protection, 269; gun violence, 263; internet connectivity, control of, 259; prohibition (alcoholic beverages), 275; regulation enforcement, 273; taxing AI, 280; unemployment insurance, 280

> 360 INDEX

United States, social media: content web scraping, 115, 191, 289 moderation, 205-9; cultural Web3, 233-34 incompetence, 192; DARPA, 259; weights, 106-7 luck, role of, 82; misinformation, 193; suicide/self-harm, 201-5 UnitedHealth, 50 Universal Basic Income (UBI), 279-80 Universal Pictures, 82 University of Chicago, 250 unpredictable events, 75-76 Upstart's model, 42 U.S. copyright law, 123, 208 U.S. Federal Trade Commission (FTC), 25, 273 U.S. National Labor Relations Board (NLRB), 279 U.S.-centric policies, 192 us-versus-it argument, 168 Utah (Live Time), 230 Uyghur population, 127

valleys, (winters), 235 vector, 119-20, 121 vector similarity, illustration of, 121 venture capitalists, 234, 240 Verizon, 259 village idiot, use of term, 170 Vinsel, Lee, 253 viral hits, 86-90 virality, 86-90 von Neumann, John, 63

Wang, Angelina, 24, 42-43 Wang, Rona, 103 WatchMojo, 208, 209 watermarks, 123 Watkins, Elizabeth, 239 Watson, Emma, 142 weather forecasting, 62-64, 65 weather prediction, 67-68

Weizenbaum, Joseph, 165 welfare checks, 203-5 WhatsApp, India, communal violence, White, E. B., 82 White, Molly, 234 White House executive order on AI, 270 Whittaker, Meredith, The Myth of Artificial Intelligence, 252 Williams, Adrienne, 146 Williams, Robert, 15 Winecoff, Amy, 239 winters (valleys), 235 Woodruff, Porcha, 15 work, AI and the future of, 276-81 World Wide Web, 233 worms (viruses), 176

X (formerly Twitter): account suspension, unaccountable, 213-14; bridging-based rankings, 220; content moderation, 180; content moderation mistake, 184; meme lottery, 87; Turkey, blocking opposing voices, 216 x.ai (calendar scheduling company), 230

You Tube: Charlie Bit My Finger, 86-87; Content ID, 206-9; content moderation mistakes, 184; copyright, exception to 1996 law, 206; copyright strikes, fake, 208-9 Your Face Belongs to Us (Hill), 16

Zeihan, Peter, 91 zero-day vulnerabilities, 175 Zuckerberg, Mark, 179, 220, 221