CONTENTS

1.	Decisions: How We Decide and Why It Matters: Human Agency and Changing the World	1
2.	Rules: The Governance of Cyberspace Offers a Cautious Tale of Hype, Hope, and Failure	18
3.	Falsities: Two Ways to Approach the Problem of Misinformation	37
4.	Bias: Why We Can't Expect AI to Solve Deep-Rooted Flaws in Human Decision-Making	59
5.	Doubt: Incomplete Information and the Problem of Irreversibility	81
6.	Principles: Guardrails Should Empower Individuals, Be Socially Anchored, and Encourage Learning	100
7.	Self-Restraint: How to Avoid the Governance Trap of Too Much Context-Awareness, or Not Enough	122
8.	Range: Four Case Studies That Illustrate the Art and Science of Making Innovative Guardrails	143
9.	Machines: Why Technology Is Neither Anathema nor a Panacea, But a Valuable Piece in the Puzzle	164
10.	Futures: How to Think About the Exercise of Power as Humans Approach a New Digital Frontier	184

viii CONTENTS

Acknowledgments	193
Notes	195
Index	219

1

DECISIONS

July 1, 2002, was a dark summer night at the German/Swiss border. Well above the clouds, a Russian Tupolev 154 airliner was cruising westward. Inside it, dozens of gifted children from Ufa, southwest of the Ural Mountains, were looking forward to a holiday in Spain. In the cockpit, highly experienced captain Alexander Gross had the controls, assisted by four colleagues. Not far away, a Boeing 757 freighter was flying northward to Brussels at the same altitude.

Noticing the converging flight trajectories, an air traffic controller for Swiss air space contacted the Tupolev crew to resolve the issue. He instructed Gross to descend and the Tupolev's crew complied.

However, both airplanes were equipped with automatic collision warning systems. Just after the air traffic controller had issued his command to descend, the collision warning systems instructed both crews to take evasive maneuvers—but it ordered the freighter to descend, and the Tupolev to climb.

Having received conflicting information from the human air traffic controller and the automated collision warning system, the Tupolev crew debated whether to continue its descent or climb instead. Their discussion was interrupted by the air traffic controller instructing them again and this time urgently to reduce its altitude, unaware that the automated system was now issuing contradictory instructions. As the crew continued on its downward trajectory—heading straight for the freighter which, following the orders of the automated system, was also

1

2 CHAPTER 1

descending—the warning system in the Tupolev more strongly commanded Gross to climb.

Collision warning systems in airplanes close to each other get in touch automatically and hash out which airplane is to climb and which to sink, to guarantee sufficient spatial separation between them as long as the system's commands are followed strictly. Hence, today standard operating procedures mandate that commands of the collision warning system must be complied with immediately, even if contradicting human air traffic controllers. But at the time, the pilots' training was not entirely clear on this matter. Forced to choose between human and machine, Gross chose to rely on the human controller. Shortly thereafter, at around 35,000 feet, the Tupolev collided at full speed with the Boeing freighter. Everyone on board both planes perished that night, high above the German city of Überlingen.¹

The accident was quickly blamed on the air traffic controller, who was overworked and with some equipment not fully functional. But there is a more fundamental issue at play. On that fateful night, the Tupolev crew faced a consequential decision: Should they trust the information coming from the human controller or the collision warning system?

True, without the air traffic controller's mistaken information to descend, the crash would not have happened. But the midair collision wasn't caused only by bad information. Gross knew he had to *choose* between good and bad information, he just was unsure which was which. Rather than asking the air traffic controller for clarification or following the warning system's advice, he *chose* to descend.

Like pilots, we too face many decisions every single day, although few of them are similarly consequential. In deciding, we rely not only on information and our own thinking. Our decision-making is also shaped by external forces, especially society, prodding, nudging, or pushing us toward a particular option, like the collision warning system. We call these *guardrails*—and that's what this book is about, from the enablers and constraints of the information we receive to rules and norms that shape how we choose among our options and how bound we are by the choices we make.

DECISIONS 3

The concept of such societal guardrails is a metaphor borrowed from the kind of physical structures you see along the sides of roads or boats. Done well, these structures offer the best of both worlds. They show you where the edge is, making it less likely that you'll step over without meaning to. But they aren't like prison walls, which make it impossible to climb over if you want. You can still go off road or take a swim if you desire. Guardrails are more about marking zones of desirable behavior rather than pushing narrowly for a single "right" choice.²

Decisional guardrails are the interface between a person's choice and the input of society. They link the individual and the collective. Decisions taken by individuals or small groups can shape the lives of many others, as the midair crash above Überlingen so horrifically exemplifies. In a world in which decision-making is largely individual, decisional guardrails are society's most direct way to influence our mutual trajectory. This book details how, collectively, we aim to alter the decisions that are being made. It is about how society governs the contexts in which individuals make decisions—a topic both powerful and ubiquitous, yet rarely understood comprehensively.

Selecting the appropriate qualities for these decision guardrails is critical. But we will argue that in our digital age we are too quick to opt for certain types of guardrails. Without much reflection, we amplify some guardrail qualities as we overemphasize the role of technology, reflecting a widespread trend for technology to increasingly govern all kinds of human decision-making. The 2002 midair collision over Überlingen seems to confirm these beliefs: If only humans follow machines, disasters are avoided.

In this book, we suggest that such a strategy is deeply flawed. This is not because technology is somehow unable or unfit to provide effective decision governance, but because the real issue is not the nature of the decision guardrails—whether they are technical or social—but the principles underlying their design. The real question is: What kind of decisions do guardrails facilitate and what decisions should they enable?

In the nine chapters that follow we examine guardrails in a variety of challenges, contexts, and cases. But our aim is not to examine every

4 CHAPTER 1

aspect or offer a detailed blueprint; we train our eye on what we think is an emerging bigger picture—a crucial red thread in appreciating the importance of designing good guardrails. Our goal is twofold: to broaden our normative horizons, so that we realize the breadth and depths of the solution space of possible guardrails; and to offer guidance that can help us craft and select guardrails that are fitting for our challenging times—to ensure not just human agency, but human progress.

Before we can fashion a solution, however, we need to better understand what's at stake and why.

Choices, Choices Everywhere

We all make decisions—hundreds, even thousands of times every day.³ Most of these decisions are trivial. We make them quickly and without much thinking. For others, often more consequential ones, we spend hours agonizing. Each decision shapes our future. The academic field of decision science is relatively young, having formally been established in the twentieth century. The quest to make good decisions, however, is as old as the human capacity to reflect on the choices we face.⁴

Relevant information is an obvious and crucial element of good decision-making. We glean insights from our social interactions with others, aided by the evolution of language. Script made it possible to preserve knowledge across time and space. Libraries, a cultural invention built on reading and writing, have served for many centuries as crucial social institutions enabling us to collect information, learn from it, and use it to make life better.⁵ The information stored and curated in these vast collections shaped decisions that led to important advances in areas as diverse as agriculture, architecture, medicine, art, manufacturing, and war. In the United States, libraries were assigned a crucial role at the birth of the nation: The Library of Congress was tasked with collecting the world's knowledge, and a nationwide system of public libraries aimed to bring this knowledge to the people. 6 The US Constitution makes clear that information is preserved and made available for a purpose, much as patents are granted not to reward the inventor, but "to promote the progress of science and useful arts."7 It recognizes that the

DECISIONS 5

role of information, in all its mediated forms, is deeply utilitarian—improving individual and societal decisions.

More recently, digital technologies have dramatically promised to lay the groundwork for better decisions by unlocking the power of computing, data, and algorithms. More than ever before, information is at the center of our daily decision-making: We consult Siri about the weather forecast, ask ChatGPT for a couple of dinner jokes, and heed Tinder's recommendations for our next date. And indeed, in the grand scheme of things digital tools have improved the conditions for decision-making, from search engines to forecasting the spread of a virus to detecting credit card fraud from subtle anomalies in transaction data.

Information we receive needs to be analyzed and evaluated. We constantly "frame" information through our mental models about how the world works, often without much conscious thought. This is what we mean when we say that we put information into perspective. This process enables us to generate and compare options. We tend to evaluate options for hugely consequential decisions more carefully, although our judgment isn't perfect—but sometimes we also fret over trivial decisions or choose bluntly without much consideration. As we ponder options, we wonder how irrevocable our actions will be. Are we bound by them, or could we reverse course if necessary?

Pop psych literature and management training courses offer a plethora of tools and tricks to help us in this process of generating and evaluating options. We are told to "think outside the box," or make a list of pros and cons. Not every such suggestion is backed up by solid research. We can't think outside the box, for instance, in the sense that we are always thinking within mental models (and decide badly if we try without them). But many suggestions may be useful in appropriate contexts.

At this point some notes of caution are in order. We are focusing here on the elements of human decision-making and how to improve that process. But we are not suggesting that all our decisions are carefully thought through. While much of our argument applies for all decision contexts, it is strongest and most valuable when we decide deliberately.

Neither are we implying that decision-making is a clean linear process, with one step followed logically after the other: collect information,

6 CHAPTER 1

analyze it using our mental models to generate decision options, compare, and choose between them. On the contrary, these elements are linked in many ways. Even deliberate decision-making is often messy and iterant. For instance, as we compare options, we may realize we missed an important dimension and must go back and gather additional information.

Nor are we suggesting that even deliberate decisions are entirely rational. Research has impressively shown that our decision-making is shaped by cognitive biases that influence our thinking. We cannot switch them off—at least not easily and at will. This realization may shatter any simplistic hope that we can achieve objective rationality in the choices we make, but it isn't fatal to the idea that the decision process is open to improvement toward better reasoning.

Decisions are important because they prepare us to take actions that shape the world. But it's not just that decisions change the world—it's that we change the world that way. Decisions are expressions of human agency—of our ability to influence the trajectory of our own existence and that of our species, even if only slightly. Human agency makes us matter. Without it, there would be no motivation to act. Agency is the source of energy that gets us out of bed in the morning to weather the storms of our daily lives.

Of course, we do not know whether we really have agency. Perhaps, from the vantage point of an omniscient objective bystander, both our actions and our sense of agency are just the results of biochemical processes over which we have no control. ¹¹ But for us, the view of the nonexistent bystander is largely irrelevant. What matters, pragmatically speaking, is what we perceive every time we select an action and take it. Consequently, in this book we embrace human agency as something that we experience as existing.

Guardrails as Governance

Decisions are the cognitive mechanisms through which we interact with the world. Much hinges on them. Understandably, society has taken a keen interest in facilitating that we decide well.

DECISIONS 7

Information is an important ingredient for good decision-making. And so, a variety of guardrails exist that shape what information is available. For instance, in the United States, corporate disclosure laws limit what a company's executives can share publicly and when. 12 Share too much information and you risk being fined, as Elon Musk found out when he tweeted about taking Tesla, a listed company, private in 2018.¹³ In other contexts, the reverse is true, and one is required to make public certain information. Pharma companies need to disclose possible side effects for the drugs they manufacture, car companies need to publish emissions and fuel efficiency figures, and the food industry needs to put nutritional labels on most of their products. 14 Sometimes, such a l'obligation d'information, as the French call it poetically, may apply to a company's clients. Insurance policies are an example. The insured is typically under a duty to disclose material facts that affect the risk to the insurer. In a similar vein, the state itself makes available a wide variety of information to help individuals make better decisions. 15 Laws are made public so that citizens can obey them, at least in democratic states. Public registers, such as for corporations or landownership, help people decide whether to engage in a business transaction.

It is not only legal rules or government policies that mandate the sharing of information. It can also be a social norm, rooted in culture and custom, such as conflict-of-interest statements in academic publications. Or it can be a practice an organization voluntarily submits to. Think, for instance, of corporate disclosure of social and environmental responsibility metrics.¹⁶

The hope behind all such interventions is that providing relevant information leads to better choices. When IKEA provides detailed instructions on how to assemble their furniture, they hope it will lead to decisions that make one's sofa bed more stable. When regulators mandate labels on food wrappers, they hope information about high calories and excessive amounts of sugar will lead people to make nutritious choices—though the chocolate bar might still be too hard to resist.

In the preceding examples, information is required in situations where a decision is imminent. In other contexts, information is meant to serve as a foundation for actions further down the road. It becomes

8 CHAPTER 1

an accountability tool with a longer shelf life. For instance, freedom of information mandates, so the theory goes (as usual, myriads of practical issues mess with the theory), enable citizens to make better decisions about the policies that affect their lives and, ultimately, give a thumbs up or down when the government is up for reelection. ¹⁷ Ralph Nader, the famous US government reform and consumer protection advocate, summarized it succinctly: "Information is the currency of democracy." ¹⁸

Beyond facilitating the flow of information, guardrails extend to the process of creating and weighing decision options. For example, numerous legal rules aim to ensure that individuals can decide without undue duress, including making extortion and coercion criminal offenses. ¹⁹ In some countries, certain particularly consequential transactions must be done before public authorities or involve testimony from experts to make sure that all parties are aware of and have considered all effects. ²⁰ Nowhere is this more evident than in the growing number of nations that have chosen to permit assisted suicide. The decision to end one's life is so grave that these societies require multiple formal steps to confirm that the decision is deliberate, free of duress, and often in the context of a terminal and painful illness. ²¹

Sometimes long-term decisions come with waiting times or "cooling-off" periods to give people ample opportunity to carefully think through their choices. ²² Being bound by a decision for a long time may have benefits—it offers stability. But we might want to think harder about whether it is the right option—and we may need more time to do so. In numerous other instances, societal guardrails explicitly enable decisions to be retracted and minds to be changed, even if that causes headaches for other parties involved. ²³

As with guardrails on information flow, guardrails on weighing options cover a spectrum from community practices to formal legal requirements. The standard operating procedures for aircraft pilots we mentioned at the start of this chapter—including whether to follow the commands of the collision warning system or the air traffic controller—are not formal law, but airlines require their flight crews to adhere to them. Similarly, emergency doctors in many hospitals must work through standard protocols of diagnosing and treating patients. It's

DECISIONS 9

not the law, but part of the organizational and professional culture—and it has been shown to be highly effective.

Such codes of conduct exist for many professions and organizations. Ever wondered how Amazon or McDonald's handles transaction complaints? They have detailed rules for how a customer service rep may decide and under what circumstances. Among merchants more generally, rules evolved over centuries that set out how they ought to behave when interacting with each other. Stemming from annual trade fairs in European cities from the thirteenth century, these rules, sometimes called "lex mercatoria," aimed to enhance trust in the market overall.²⁴

A far more subtle shaping of individual decision processes has become popular lately in some policy circles. Called "nudging," the idea is to delicately prompt people to choose the option that will be most beneficial for them. For example, when it is judged that not enough individuals opt into a retirement savings plan, one could make participation the default and require those who do not want to partake to actively opt out instead. Advocates tout nudging as less limiting than more outright restrictions, but skeptics point out that nudges are opaque, creating an illusion of choice while manipulating the decision process. Each of the process of

Similar techniques can be used to shape decisions in ways that further the interests of people other than the decision-maker. Ads and salespeople use a wide variety of cognitive tricks to influence transaction decisions. Even the layout of supermarkets is carefully designed to affect our purchasing choices. Deep-rooted social and cultural practices can be deliberately repurposed to shape our decisions. In the early years of eBay, sellers often rated buyers highly *before* a transaction had been completed. That didn't make sense. Why should you rate somebody before you know whether she did as promised? Researchers took a closer look and discovered that such a premature positive rating was perceived by the buyer as a gift, which gave rise to a social expectation to reciprocate. Those who quickly rated the other side in positive terms got more favorable ratings in return, which somewhat divorced ratings from the underlying transaction and prompted eBay to change its rating system.

10 CHAPTER 1

So far, we have drawn a distinction between measures that shape the information we receive and measures that influence how we evaluate decision options. The distinction is artificial, in the sense that all measures that shape our decision processes involve information—otherwise they would not be able to reach into our mind. Airlines' standard operating procedures shape how pilots weigh their options, but they are also information that pilots read and digest. When a nudge shifts a decision default, it's also information about how easy or hard it is to decide on a particular option. However, we find the distinction between "informational" and "decisional" guardrails useful because it helps us comprehend the wide spectrum of possibilities.

As will be clear from the examples above, by guardrails we mean more than a simple norm or rule. Guardrails often include processes and institutions, mechanisms and tools, even a "culture" or "way of thinking." For instance, emergency doctors have internalized checklists, while standard operating procedures can take on material form in safety mechanisms in factory machinery. Programmers at large software companies live by a "software development life cycle," a combination of rules, processes, and organizational structures to help ensure good coding. ³⁰ It is the "system" around a naked norm—the processes and institutions—that makes guardrails work. Hence, when we write here about guardrails we see beyond single rules and include the reality around them that makes them work (or not).

Because our notion of guardrails isn't limited to formal legal rules and because we include the structures around them, we see them in a very wide variety of contexts and circumstances. Dynamics like globalization have, some scholars maintain, proliferated the types and kinds of guardrails, leading to a pluralization of regulation.³¹ Others, like Gillian Hadfield, agree—and turn the analysis into a prescription, suggesting we need to think more in terms of markets of rules than a hierarchy of them.³² Whatever the concrete causes and consequences, what matters in our context is simply that guardrails shaping our decisions are plentiful and diverse. But if decision-making is the expression of an individual's volition, why are others—communities, society—so interested in shaping individual choices?

DECISIONS 11

The Social and Externalities

The obvious answer is that as social beings, we care for each other. Helping each other is something we practice right from early childhood, so why should we not want to help each other to make good decisions? A friend, colleague, or a complete stranger may benefit from measures that improve their individual decision-making today—but we may be the lucky recipients of guidance tomorrow.

Anthropologists offer another compelling argument. Humans have made stunning progress over the past few millennia—compared to other species, but also to earlier phases of human existence. This cannot be explained by biological evolution, as the cogwheels of natural selection do not operate fast enough. Instead, what has propelled us forward so dramatically is some form of *cultural* evolution that involves learning things from each other, rather than having to learn everything for ourselves.³³ It's a marvelous cognitive shortcut to discovery: Insights can be passed on. We can stand on the shoulders of those who came before us. The key is our ability to learn abstractly, to let our minds wander instead of our bodies. When it comes to decision-making, too, communities want to ensure that good insights spread. We are eager to share suitable guardrails and are open to accept them—at least to an extent.

Economists put forward a related but distinct reason for societal guardrails. When people make decisions that affect other people, economists call those effects externalities. Implementing guardrails can serve a utilitarian purpose, as shaping an individual's decision influences the externalities the decision causes. For example, in 2015, the US Environmental Protection Agency (EPA) discovered that Volkswagen, one of the world's largest car manufacturers, had illegally deployed software in more than ten million of its cars to deceive emissions tests—thereby evading a requirement to provide truthful information that will help car buyers make good decisions. Top managers and engineers at the car company had known about the illegal scheme for years. ³⁴ As a result, millions of consumers bought cars erroneously certified as green and powerful, which caused huge amounts of unhealthy emissions. When

12 CHAPTER 1

the deception became clear, millions of affected cars lost much of their residual value overnight.

Externalities can be positive as well as negative, of course. The decision of a well-known coffee chain to open a shop in a troubled neighborhood can be seen as a signal of confidence in the neighborhood's future, attracting others to invest and creating new opportunities for people nearby.

Decisions can have consequences that impact groups and institutions as well as other individuals. This idea is illustrated by the textbook example of used car sales. A car's history—such as whether it has been in accidents or has serious mechanical problems—is not always evident from looking at it. Absent any requirement to disclose such information, buyers tend to distrust used cars. They bid less than they would if they knew the car was good, because they factor in the risk that the car may be a "lemon." This is unfortunate for the honest seller, who will not get the car's actual worth from a sale.

There is a bigger and more pernicious consequence, though. Discouraged by not being able to sell their good cars for a fair price, honest sellers exit the market. As economists have shown, this leads to a vicious cycle: As lemons account for more of the market, buyers become even more reluctant to transact. This makes the *market* ineffective. The societal intervention in many nations in response is to require sellers to disclose whether their car had previously been in an accident. This not only helps buyers make the right decision, it helps honest sellers to find buyers—which increases the average quality of cars on offer in the market, and enables the market to do what markets should: help allocate a scarce resource.

Because collectively we benefit from better decisions, for society it makes sense to establish guardrails to inform and affect decision-making. By influencing individual decisions, guardrails enable society to chart a middle path between two extremes: full individualism, unencumbered by collective needs, or complete control through the collective without regard for individual preferences. Instead of a choice between Ayn Rand's *Atlas Shrugged* and George Orwell's *1984*, good societal guardrails offer the best of both worlds—exercising societal control without negating individual volition.

DECISIONS 13

Good guardrails are a sweet spot that is challenging to find. In the abstract, they effectively guide appropriate individual decision processes—but in concrete contexts, defining effective and appropriate ones is a difficult if worthy challenge. Guardrails are not only, to quote political scientist Friedrich Kratochwil, "guidance devices" shaping individual decisions, but "also means which allow people to pursue goals, share meanings, communicate with each other, criticize assertions and justify actions." They signify that individuals are being taken seriously not only as decision-makers, but as members of the society they live in. There are a lot of moving parts to keep in mind when crafting a good guardrail. But rather than tackling this challenge head on, in recent years we have become sidetracked by technology.

The Technological Digression

Humans have used technical tools to aid their decision-making for centuries, but digital technologies now promise to be an unprecedented turbo for improving our decision-making—unlocking information bottlenecks and offering humans comprehensive access to knowledge.

Take only generative artificial intelligence systems like GPT (Generative Pre-trained Transformer). Trained by ingesting more than half a trillion almost entirely human-written words from millions of digitized books as well as billions of web pages, GPT is built on the collective knowledge and experience of humanity (or at least a significant slice of it).³⁸ It is being used to retrieve decision-relevant facts and information, but also to provide a wide variety of known decision options and to even offer decision recommendations.

However, these technologies have given rise to new and urgent questions about who controls digital information flows and the algorithms that power them. Data-driven machine learning models, such as GPT, are black boxes; we can interact with them but not peek into them easily. It's not just cutting-edge AI that is incomprehensible, however. Infamously, it is said that not even Google's own engineers can fully understand the inner workings of their search engine, upon which so many of us rely to inform our daily decisions. ³⁹ Nor do we know what exactly

14 CHAPTER 1

determines the news feed on our favorite social media channels, or what shapes which ads we get to see when browsing the web.

Faced with questions about the complexity of such systems, the response often is that we need *more technology, not less*. AI is touted as capable of making better-than-human decisions. More data, improved algorithms, and more computing power promise to govern decision-making processes of the future. From smart contracts to autonomous vehicles, decisions are increasingly prepared, performed, and executed within technical systems.

In part, this focus on technology is also a consequence of the rise of "Big Tech"—companies operating digital platforms that have accrued enormous power to shape information flows for billions of users. However, it also reflects a process that has been going on for many decades, even centuries. ⁴⁰ To describe this process, historian Lorraine Daston differentiates between "thick" and "thin" rules. Thick rules require interpretation and social acceptance. This means they may not be perfectly enforced, but their flexibility often makes them effective. In contrast, thin rules are stepwise instructions that are set and fixed. It's hard to derail them, so they can be relied on—but they can't be adapted easily.

Daston argues that in Western societies thin rules have risen, while thick rules have declined. She points to the rise of the administrative state and detailed and comprehensive regulatory rules that try to cover all eventualities in advance. This brings more predictability and lowers risk; but it also means that there is less room for discretion, change, and flexibility. Algorithms are an example of thin rules. Their increasingly extensive use as guardrails can be seen as a continuation of a process that started long ago.

In this book we challenge the mantra that more technology is the best answer to problems of human decision-making. Of course, we acknowledge that technology has the capacity to empower individuals, institutions, and society at large to make better decisions. ⁴¹ We also recognize that technology is never neutral: What technology is chosen has consequences for what can be achieved with it and how. ⁴² And we agree that the link between technologies and commercial control ought to be scrutinized: Opaque values baked into crucial technical bottlenecks

DECISIONS 15

of global information flows that influence billions of individual decisions need our critical questioning.⁴³

But our concern is more fundamental. We argue that the focus on technology is distracting because it shifts our attention to a discussion over operational mechanisms and their implications when instead we should be engaging in a normative debate about the right qualities of guardrails.

Qualities for Times of Uncertainty

Every shift of our focus comes at a cost. When light is shone on one feature, others remain in the dark—understudied and overlooked. And so, by focusing on technology, we lose sight of what we suggest matters most in our times: defining the qualities and features of our society's decision governance.

Our starting point is the proposition that our world is becoming more volatile and uncertain. Challenges as diverse as social justice, public health, geopolitical disorder, and climate change will persist and deepen in the decades to come. The frameworks we put in place today to guide our decisions must be able to strategically embrace uncertainty. But the technologies that promise to improve decision-making regularly seek to negate flexibility and uncertainty, as we will discuss in chapters 3, 4, and 5.

So what are the alternatives? How can we create and employ guardrails that support and guide our decision-making in a world marked by increased uncertainty?

We argue that we need to understand our situation as one that requires less *technical* innovation than *social* innovation. We need to build on existing processes and institutions. The real challenge lies less in the concrete mechanics of guardrails than in getting the foundations right. We know this from our everyday practices. Before a driver revs up her engine, she needs to clarify where she wants to go, and what aspects of the journey—speed, safety, cost—she most cares about. We, too, must first clarify not just what our goals are, but also what qualities we want to have embedded in the mechanisms we employ to reach these goals.

16 CHAPTER 1

We need to choose the decision qualities we want our guardrails to further and facilitate. This necessitates analysis and critique, but also normative thinking, both about society's role in providing these guardrails and what environment for individual decision-making we envision. We map out a suitable process and develop three concrete design rules for good guardrails in chapter 6. We then add an important, perhaps crucial, constraint to guardrail design in chapter 7.

As we do this, we will realize that we already have solid foundations on which we can build. We recognize these qualities in some of the governance mechanisms we already employ—with positive results. And we will see that the space for governance mechanisms to incorporate some (or all) of these desired qualities is far larger than we initially might have believed. A fresh but detailed look at the qualities inherent in various kinds of guardrails can help us see a broader spectrum of possibilities. By combining mechanisms and institutions, we can establish the innovative governance framework we need. In chapter 8 we'll map this diverse governance landscape in greater detail.

Implementing this governance framework may involve the use of technology, but only to the extent that it advances our goals and reflects the qualities we seek; we show how in greater detail in chapter 9. To foreshadow, we need to be less impressed by superfast bits traversing cutting-edge hardware than by existing social mechanisms that have proven their use. Rather than supplanting existing governance setups, technology should support them as a tool.

Widening the Aperture

Decisions prepare us to take actions. Through our actions we change reality. Humans aren't the strongest or fastest species. We may have mastered arithmetic, but computers calculate faster than we can. We may be excellent at recognizing shapes and patterns, but AI turns out to be even better. So, what's left for us?

As humans we believe in agency—in our ability to choose and shape the lives we live. Steve Jobs referred to it as the desire to "leave a little dent in the universe."⁴⁴ But the desire to make a difference manifests

DECISIONS 17

itself not just on the individual level. As a society, even as a species, we want to effect change. When Neil Armstrong stepped on the moon, he said it was a "giant leap for mankind" because it showed how humanity could accomplish that little dent in the universe.

Societal guardrails to individual decision-making are where the collective and the individual meet. Considering what decision qualities they should enable means asking what is right for both the "I" and the "we." We cannot define society's goals without conceptualizing what we want individuals to aim for. Through guardrails, society may express itself by injecting its values into individual decision-making.

Thinking normatively about societal guardrails also entails pondering the role of the individual in society. Four decades ago, US constitutional scholar Kenneth Karst suggested the metaphor of "equal citizenship," capturing an individual's equal agency as part of a greater compact. ⁴⁵ Around the same time, but across the Atlantic, the German Constitutional Court opined eloquently about "informational self-determination" as an "I" that is always contextualized and anchored in a "we." ⁴⁶ Harvard's Human Flourishing Program emphasizes the human need to evolve along five dimensions, from the highly individual to the deeply collective. ⁴⁷ The message of these three and many others is clear: We are individuals, but we also are a part of something larger.

So as we write about the qualities of guardrails in the chapters that follow, we are not only opining about society and its role. We are also writing normatively about the individual: the place she ought to occupy, the values she ought to cherish, and the goals she ought to attain. Guardrails are, to paraphrase sociologist Anthony Giddens, "social practices"—structural mechanisms that reconfigure and reshape society. ⁴⁸ If our initial focus may seem narrow, the implications are far bigger. Because through the decisions we make, we become not only agents of our destiny, but fellows of our society.

INDEX

ABS systems, 71 ambiguity, 97-98 Amundsen, Roald, 100 Adams, John, 137 advertising, 63 Anderson, Benedict, 107 agency: decisions as expression of, 6; guard-Apple, 32 rails as empowerment of, 108-10; human apprenticeships, 64 condition and, 113-14, 190; innovation as Aquinas, Thomas, 135 instance of, 74-76; learning linked to, Arendt, Hannah, 42 113-14; progress resulting from, 79; Aristotle, 135 societal, 16-17; technology as defense Armstrong, Neil, 17 assisted suicide, 8 against, 66-67 AI (artificial intelligence): bias in, 177-82; autonomous vehicles, 81-84 criticisms of, 71-73; decision-making AWS (Amazon subsidiary), 179 role of, 13-14; information filters, 51-52; limitations and drawbacks, 73-77, Bambauer, Derek, 71 79-80, 102; machine learning, 68-77, Beatles, 158-59 Benedict, Saint, 185-87 129, 177-79 AI4ALL, 181–82 Benkler, Yochai, 214n30 Airbnb, 87 best practices, 65, 149-51, 162 Airbus, 66 bias: AI and, 180-82; algorithmic, 175-81, air traffic controllers, 1-2 215n19; cognitive, 6, 62-64; in lending, Algorithmic Justice League, 181 173-75, 215n19; machine learning and, algorithms: bias in, 175-81, 215n19; content 72-73; types of, 62 selection conducted by, 50; as informa-Bible, 41 tion filters, 48, 53, 58; inscrutability of, 13; Biden, Joe, 176 lending decisions made by, 164-65, 173-75; Big Bang Theory, 65 misinformation distributed by, 40; profit black boxes, machine-learning models as, motive behind, 127; promise of, 14; recom-13, 71-73, 178-79 mendation function conducted by, 64; blockchain, 91-93 rigidness of, 14; risk calculation con-Brexit, 40, 107 ducted by, 85; underwriting conducted by, Britannica (encyclopedia), 159-61 Buolamwini, Joy, 181 215n19 AlphaGo, 69, 76 Burkina Faso, 151-54 Amazon, 129 Bush, George W., 167

220 INDEX

Cairncross, Francis, 29 Cambridge Analytica, 40 Canada, 48 cap-and-trade system, 109 capitalism, 34-35, 119 Carnegie Mellon University, 149 Cassidy, Joe, 37-38 Catholic Church, 41 censorship, 53 CERT. See Computer Emergency Response checklists, 10, 65-66, 68 chemical weapons, 37-38 China, 31, 40, 132, 155-56, 158 Christianity, 41, 65 Citron, Danielle Keats, 44 Clinton, Bill, 136 code: as law, 30-33, 51-52; smart contracts and, 91-94. See also dry code; software; wet code codes of conduct, 9 Cohen, Julie, 34-35 collision warning systems, 1-2 CompuServe Germany, 126 Computer Emergency Response Team (CERT), 147-51 Conficker (malware), 150 confirmation bias, 62-63 consistency, 105-7 consumer protection, 128-29 Content ID (YouTube), 131 context: for guardrails' use, 123-25; information affected by, 28-29, 48-49, 53, 125-27, 131; strategies for determining, 127-33 contracts (traditional), 91-92, 94-97. See also smart contracts copyright, 22–23, 32, 46, 111, 130, 166–72 COVID-19, 39, 53-54, 187-88 Creative Commons, 170-73, 176, 178-80 credit scoring, 174-75 cultural evolution, 11 Curtis, Pavel, 18 cybernetics, 140

cybersecurity, 146–51 cyberspace: communication protocols in, 116–18; copyright in, 168; dispute resolution in, 89–90; distance/time not a factor in, 29, 49–50, 141; distinct features of, 21–23, 25, 31–32; governance of, 21–36, 46, 51–52, 116–18; private vs. public control of, 32–33, 51–52; problems in, 33–35; and Russian war against Ukraine, 146–48. See also digital platforms; metaverse

D'Agostino, Abbey, 59–61, 76, 77 DAO (decentralized autonomous organization), 93-94 Daston, Lorraine, 14 data aggregation, 86-88 Davies, Alex, 69 decisions: case example of, 1-2, 59-60; complicating factors in, 6, 62-64; conformity in, 74-79; constraints on, 66; contextual factors in, 123-25; empowerment for, 108-10; evaluation of, 60-62, 77-79; guardrails' role in, 2-3, 8-10, 64-67; guides, aids, and strategies for, 5, 64-65; individuals' role in, 3–6, 17; information's role in, 4–5, 7–8, 10, 81, 83–84, 123; machine learning and, 70-77; processes, 5-6, 60-64; rationality of, 6, 63; rules' and law's role in, 65-67; self-restraint in, 135-42; society's role in, 3, 11-13, 15-17, 65, 67–68; technology's role in, 3, 13–15; time factors in, 61; in uncertain conditions, 15-16; variability in, 77-80, 144 DeepMind, 69 Defense Advanced Research Projects Agency (DARPA), 149 design principles, for guardrails and solution spaces, 108-16, 120-21, 125, 144-45, 151, 155, 162-63, 172, 179, 183, 187, 189 Digital Millennium Copyright Act (DMCA), 130-31, 168 digital platforms, regulation of content on, 45-49, 51-54. See also social media

INDEX 221

discrimination, 164–65, 173–76, 215119.

See also bias
disinformation. See misinformation/
disinformation
dispute resolution, 89–90
DMCA. See Digital Millennium Copyright
Act
dry code, 91–95

Easterbrook, Frank, 20-21 eBay, 9, 26, 32, 86-90, 92, 99 The Economist, 126 empowerment, as design principle, 108-10, 113, 114, 120, 121, 145, 155, 172, 179 enclosure, 24 end-to-end principle, 31-32 Equal Credit Opportunity Act, 174 espionage, 37-38 ethical dilemmas, 82-84, 156-58 European Court of Human Rights, 44 European Court of Justice, 49, 107, 124 European Union (EU), 27, 29, 43–44, 107, 128-29, 140; Code of Practice on Disinformation, 130; Copyright Directive, 130 Evans, Teddy, 101 experimentation, 74-78. See also innovation externalities, 11-12

Facebook, 22, 39–40, 43, 47–48, 53, 126–27, 130. See also Meta
Facebook Papers, 53
Fairfield, Joshua, 201144
Fair Housing Act, 174
fair use, 131, 167
fake news, 49–52
Fannie Mae, 174
Federal Bureau of Investigation (FBI), 37–38
federalism, 141
feedback forums, 86–89
FICO scores, 174–75, 215119
Fisher, William, III, 211120
flexible guardrails, 14, 33, 68, 79–80, 95–99, 103–4, 140, 143–46, 162, 185–86, 189–90

food labels, 128
FORVM (journal), 44
Frankel, Tamar, 116–18
Franklin, Benjamin, 42
Freddie Mac, 174
free speech, 42–45. See also hate speech
Frischmann, Brett, 201144
Frug, Gerald, 24

Gawande, Atul, 65 generativity, 115 German Constitutional Court, 17 Ghana, 151-54 Giddens, Anthony, 17, 112 Gitlow v. New York, 112 Giuliani, Rudy, 39 Glawischnig, Eva, 43-45 Global Network Initiative (GNI), 157-58 Global Online Freedom Act, 156 Global Water Partnership, 153 GNI. See Global Network Initiative Go, 69 Goldsmith, Jack, 27 Google, 13, 32, 46, 49, 69, 124, 126-27, 130, 169, 179, 198n39 Google Spain judgment, 124 governance. See guardrails; rules/law GPT (Generative Pre-trained Transformer), 13, 70

Gross, Alexander, 1–2
guardrails: case examples of, 18–19, 36,
146–62, 166–73, 185; catch-all, 132–33;
complexity of, 103–4; contextual factors,
123–25; critique and modification of,
111–15, 118, 121, 134–35, 144, 145, 146, 157,
188 (see also design principles, this entry);
cyberspace, 21–36; decision-making role
of, 2–3, 8–10, 64–67; defined, 2–3, 103–4,
143; design principles, 108–16, 121, 125,
144–45, 151, 155, 162–63, 172, 179 (see also
critique and modification of, this entry);
empowerment through, 108–10; flexible,

Great Firewall of China, 31

222 INDEX

guardrails (continued)

14, 33, 68, 79–80, 95–99, 103–4, 140, 143–46, 162, 185–86, 189–90; formal and informal, 8, 10, 112, 118; for information availability and integrity, 7–8, 84–86; learning encouraged by, 113–16, 118–19; legal, 95–99, 110–12, 119, 128, 139–40, 142, 144; on lying, 41; power issues regarding, 188–90; self-restraint in, 125, 135–42, 145; social, 3, 8, 11–12, 15–17, 64, 79, 84–86, 97–99, 101–2, 107–21, 172, 176, 178–79, 182, 186, 190; societal benefits of, 11–13, 110–11; spectrum of, 8–10, 104, 141–42; technical, 36, 80, 102–3, 112–13, 121, 164–83, 188. See also rules/law

Guthrie, Woody, 166-67, 172

Haakjöringsköd, 94–95 Hadfield, Gillian, 10 Hamblin, Nikki, 59–60, 76, 77 Hamilton, Alexander, 208n10 hate speech, 43–46, 49, 51 Henrich, Joseph, 114 *Hidden Figures*, 78 Hildebrandt, Mireille, 98 Hinduism, 41, 65 Hoffman, Reid, 92 Holmes, Oliver Wendell, 111–12 Human Flourishing Program, 17 hydroxychloroquine, 39

ICANN. See Internet Corporation for Assigned Names and Numbers information: aggregation of data, 86–88; availability and integrity of, 7–8; capitalism and, 34–35; contextual factors for, 28–29, 48–49, 53, 125–27, 131; decision-making role of, 4–5, 7–8, 10, 81, 83–84, 123; digital storage and delivery, 5; evaluation through social debate, 55–57; filters, 51–54, 57, 130; free speech, 42–45; governance of, 21–36; guardrails for, 7–8, 84–86; insufficient, 83–90, 93–94; interpretation

of, 53-57; lobbyists' use of, 28; predictive value of, 21, 67-68, 73-77, 86-88, 99, 103; reputation systems, 86-89. See also misinformation/disinformation information gaps, 83-90, 93-94 innovation: AI limitations, 74-76; in governance, 16; human-led, 74-75; social, 15. See also experimentation Instagram, 22 insurance, 85 Integrated Water Resources Management (IWRM), 153-54 intellectual property rights. See copyright intentions, 95 International Union for Conservation of Nature, 153-54 Internet. See cyberspace Internet Corporation for Assigned Names and Numbers (ICANN), 117-18, 134 Internet Engineering Task Force (IETF), 26, 116 interpretation of information, 53-54 Iran, 40 Islam, 41, 65, 106 IWRM. See Integrated Water Resources Management

Jackson, Mary, 78 Jefferson, Thomas, 42, 137 JibJab, 167 Jobs, Steve, 16, 18, 78, 170 Johnston, David, 26–27 Judaism, 65, 106

Kahneman, Daniel, 61
Kalimantan, 182–83
Kamloops Indian Residential School,
British Columbia, 48
Kant, Immanuel, 135, 136
Kaprun train disaster, 122–25
Karst, Kenneth, 17
Katsh, Ethan, 89–90
Kennedy, John F., 78

INDEX 223

Kerry, John, 167 solutions and filters for, 38, 51-54, 130; Kratochwil, Friedrich, 13 timing issues, 45-46, 50, 53; Wikipedia, 160. See also hate speech LambdaMOO, 18-21, 36 Mitchell, William, 21 law. See rules/law MIT Media Lab, 82, 181 learning, as design principle, 113-16, 118-21, moderation. See self-restraint 145, 151, 155, 161, 172, 179, 187–89. See also monastic rules, 185-87 machine learning; schooling monist law, 26 Morris worm, 149 legal pluralism, 106-7 Lemley, Mark, 207n26, 211n30 mortgage loans, 164-65, 173-76, 215n19 lending decisions. See mortgage loans Mueller, John, 198n39 Lessig, Lawrence, 23-25, 28, 32-35, 51, 111, Mueller, Milton, 209n26 Musk, Elon, 7, 78 168-70 lex mercatoria, 9 MySpace, 32 libel, 44, 45 Nader, Ralph, 8 libraries, 4 Library of Congress, 4 National Enquirer, 50 lobbying, 28 Navalny, Alexei, 38 loss aversion bias, 62 netiquette, 26 Luther, Martin, 41 Neuberger, Anne, 146-47, 151 New York Times, 129 lying, 41-42 norms. See social norms North Korea, 40, 136 machine learning, 68-77, 129, 177-79 Madison, James, 137-38 Novichok, 38 NSI (company), 116 Marbury, William, 137–38 markets, 79 nudging, 9, 65, 75 Marshall, John, 138-39 Oberschlick, Gerhard, 44-45 Marx, Karl, 34 McAfee Cyber Threat Alliance, 150 ODR. See online dispute resolution Omidyar, Pierre, 86 McDaniels, Crystal Marie and Eskias, 164-65, online dispute resolution (ODR), 89-90, 99 174 mental models, 61-62 online piracy, 130-31 opinion, 42-43 Meta, 52, 129, 130; ThreatExchange, 150. See also Facebook optimization problems, 133-34, 142 metaverse, 18, 184-85, 187, 191 originalism, 98 Orwell, George, 1984, 12, 33 Microsoft, 32, 129, 130, 179 Milton, John, 135 Palo Alto Research Center (PARC), 18 Minow, Martha, 209n25 misinformation/disinformation: case ex-Pasquale, Frank, 71 ample of, 37-38; dangers of, 39; fake news, past, as predictive of the future, 67-68, 49-52; free speech and, 42-43; interpre-73-77, 86-88, 99 tation of, 53-57; and lying, 41-42; social Perimetr system, 66-67 media and, 39-41, 50; technological Perry, William, 136

224 INDEX

Petrov, Stanislav, 135-36 Pinker, Steven, 114 piracy. See online piracy Pirie, Fernanda, 190 Pistor, Katharina, 98, 119 platforms. See digital platforms Plato, 135 Polanyi, Karl, 34 Popper, Karl, 42 Post, David, 27 Postel, Jon, 116-17 power, in formation and implementation of guardrails, 188–90 predictive value of information, 21, 67-68, 73-77, 86-88, 99, 103 Projet d'Amélioration de la Gouvernance de l'Eau dans le bassin de la Volta (PAGEV), 153-54 public domain, 166-71

Rand, Ayn, Atlas Shrugged, 12 Rawlings, J. J., 152 recency bias, 62-63 relational contracts theory, 96 religion, 41, 65 reputation systems, 86-89 right to be forgotten, 49, 124 Risch, Michael, 71 risk management, 85 Roberts, Sarah T., 47 robophobia, 81 rule of law, 139-40 Rule of St. Benedict, 185-87 rules/law: AI and, 68; code as, 30-33, 51-52; consistency of, 105-7; contracts, 94-97; and cyberspace, 21-36, 46, 116-18; decisionmaking role of, 65-67; determination of context through, 128–29; digital platforms and, 45-49; elasticity of, 95-99, 110-12, 115, 144; and free speech,

43-45; as guardrails, 95-99, 110-12, 119,

128, 139-40, 142, 144; and hate speech,

42-46; in information-based contexts,

21–36; learning and, 115; local vs. general, 25–29, 46–47, 105–7; and misinformation, 51–52; monastic, 185–87; progress resulting from, 79; proliferation of, 128–29; self-restraint in, 137–42; for take-down requests, 46–47; thick vs. thin, 14, 185–86, 189. *See also* guardrails
Russia, 40, 146–48. *See also* Soviet Union

Samarajiva, Rohan, 21 schooling, 64. See also learning Schwartz, Paul, 201140 scientific method, 42 Scott, Robert Falcon, 100-101 Second Life, 22 self-driving cars. See autonomous vehicles self-restraint, as design principle, 125, 135-42, 145, 161-62, 172, 182, 189 Selinger, Evan, 201144 Shakespeare, William, 173 shark meat, 94-95 Shi Tao Shi, 132-33, 155-57 Skelton, Reginald, 101 SMART cars, 108-9 smart contracts, 90-94, 99, 102 Smith, Brad, 129 social anchoring, as design principle, 110-13, 145, 155, 161, 162-63, 170, 172, 179, 187, 190 social media: and free speech, 43-45; and hate speech, 42-46; misinformation distributed by, 39-41, 50; as news source, 40 social norms: conflicting, 105; in cyberspace,

67–68, 101, 112 society: benefits of guardrails for, 11–13, 17; decision-making role of, 3, 11–13, 15–17, 65; evaluation of information in, 55–57; guardrails established by, 3, 8, 11–12, 15–17, 64, 79, 84–86, 97–99, 101–2, 107–21, 172, 176, 178–79, 182, 186, 190; individuals' role

19, 26; as decisional guardrails, 65; efficacy of, 120; fragility of, 33; for information

governance, 7, 25, 50, 57; law vs., 25-27,

29, 66, 105, 112; persistence/durability of,

INDEX 225

in, 17; learning in context of, 114-15; Tesla, 7 technology's place in, 15-16, 112-13, 165-66, textualism, 98 172-73, 175-83 theory of incomplete contracts, 96 software: dangers of, 33-34; as law, 30-33, thick rules, 14, 185-86, 189 51-52; private vs. government control thin rules, 14, 185-86 over, 32-33, 51-54. See also code Tiananmen Square, China, 132, 156 solution spaces: design principles, 120-21, TikTok, 64, 130 Tk'emlúps te Secwépemc First Nation, 48 125, 144-45, 151, 155, 162-63, 172, 179; as toading, 36 reservoir for governance decisions and Tomasello, Michael, 114 implementation, 4, 104-5, 116, 145, 155, 160, 187 traveling salesman problem, 134 Somm, Felix, 126 trolley problem, 82-84 SOPs. See standard operating procedures Trump, Donald, 39, 41 truth, 41-42, 55 South Korea, 136 Turing, Alan, 30 South Pole expeditions, 100-101 Soviet Union, 37-38, 66, 135-36. See also Turing machine, 112 Russia Twitter, 39-40 spies. See espionage SquareTrade, 90 Uber, 87 standard operating procedures (SOPs), 1, 8, Ukraine, 146–49 10, 65, 68, 75, 78 uncertainty: contracts in conditions of, 96; Stephenson, Neal, 18, 184 in cyberspace, 21; decision-making in surveillance capitalism, 34 conditions of, 15, 81-82, 99, 103; effect of Susskind, Jamie, 201144 local vs. general rules on, 25, 27; Rule of St. Benedict and, 185; technological Szabo, Nick, 91 attempts to mitigate, 15, 81, 85-86 take-down requests, 46-47, 49 underwriting algorithms, 215n19 Tamanaha, Brian, 106 unexpected circumstances, 89-90, 99, technology: decision-making role of, 3, 13-15; 187-88 University of Massachusetts, 89-90 determination of context through, 129-31; guardrails constituted by, 36, 80, 102-3, University of Oxford, 178-79 112-13, 121, 164-83, 188; information filters, US Constitution, 98, 138 51-54, 57, 130; information use and storage, US Department of Commerce, 117 5; limitations and drawbacks, 14-15, 33, US Department of Defense, 149 52-54, 57-58, 102-3, 120-21, 186; promise US Department of Homeland Security, 149 of, 15, 81, 85-86, 101-2, 175, 186; for used car sales, 12 redressing technological wrongs, 178-81; US Environmental Protection Agency social embeddedness of, 15-16, 112-13, (EPA), 11 US Food and Drug Administration (FDA), 165-66, 172-73, 175-83. See also software Telecom Act, 27 128 temperance. See self-restraint US National Institute of Standards and

Technology (NIST), 180-81

US National Security Agency (NSA), 147

Terra Nova Expedition, 100

terrorism, 130

226 INDEX

US Securities and Exchange Commission

(SEC), 109

US Supreme Court, 98, 107, 111–12,

137-39

utilitarianism, 83

vending machines, 91

virtual space. See cyberspace

volition. *See* agency Volkswagen, 11

Volta River Basin, 151–54

Wales, Jimbo, 160, 162

Wall Street Journal, 158–59 wampum belts, 48

Water Governance Project, 154

Waymo, 73–74 Werbach, Kevin, 92 wet code, 91–97 whale meat, 94–95 Wikipedia, 158–62, 170

will. See agency

Windhager, Maria, 43-45

Wise, David, 38

World of Warcraft, 22

Wu, Tim, 27

Xerox, 18

Yahoo! 155-58

YouTube, 32, 46, 130–31, 169

Zittrain, Jonathan, 115

Zuboff, Shoshana, 34-35, 119

Zuckerberg, Mark, 18 Zurich Insurance, 129