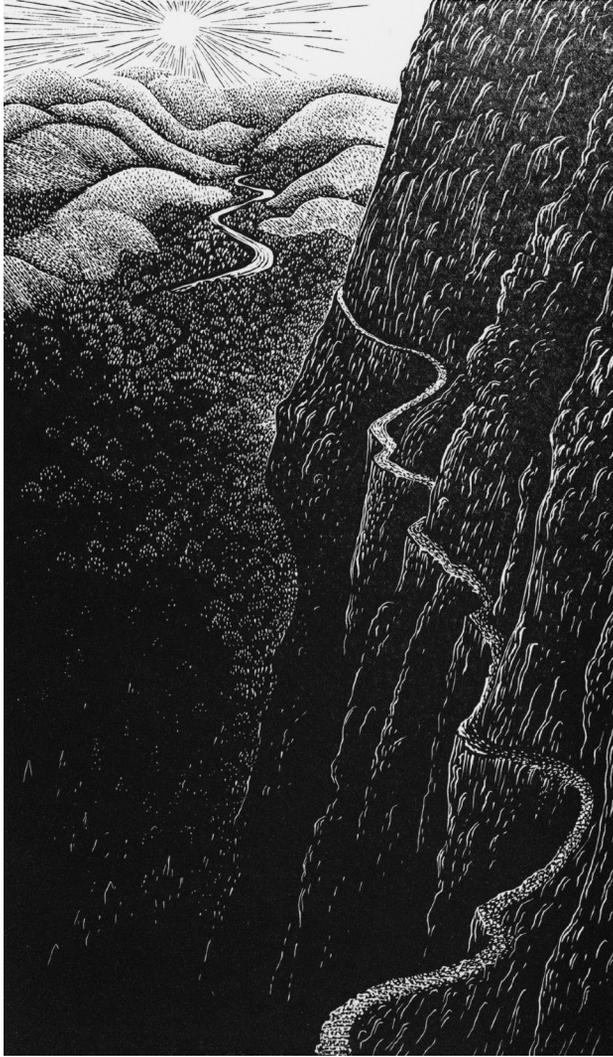


The Precipice
Existential Risk
and the Future
of Humanity
Toby Ord

B L O O M S B U R Y

THE PRECIPICE



THE PRECIPICE

*Existential Risk and
the Future of Humanity*

TOBY ORD

B L O O M S B U R Y P U B L I S H I N G
LONDON • OXFORD • NEW YORK • NEW DELHI • SYDNEY

BLOOMSBURY PUBLISHING
Bloomsbury Publishing Plc
50 Bedford Square, London, WC1B 3DP, UK

BLOOMSBURY, BLOOMSBURY PUBLISHING and the Diana logo are
trademarks of Bloomsbury Publishing Plc

First published in Great Britain 2020

Copyright © Toby Ord, 2020

Frontispiece illustration © Hilary Paynter, 2020

Toby Ord has asserted his right under the Copyright, Designs and Patents
Act, 1988, to be identified as Author of this work

Extract from *Pale Blue Dot* copyright © 1994 Carl Sagan. Originally published in *Pale
Blue Dot* by Random House. Reprinted with permission from Democritus Properties,
LLC. All rights reserved this material cannot be further circulated without written
permission of Democritus Properties, LLC.

Extract from *A Choice of Catastrophes: The Disasters That Threaten Our World* by
Isaac Asimov. Copyright © 1979 by Isaac Asimov. Reprinted with the permission of
Simon & Schuster, Inc. All rights reserved.

All rights reserved. No part of this publication may be reproduced or transmitted
in any form or by any means, electronic or mechanical, including photocopying,
recording, or any information storage or retrieval system, without prior permission
in writing from the publishers

Bloomsbury Publishing Plc does not have any control over, or responsibility for, any
third-party websites referred to or in this book. All internet addresses given in this
book were correct at the time of going to press. The author and publisher regret
any inconvenience caused if addresses have changed or sites have ceased to exist,
but can accept no responsibility for any such changes

A catalogue record for this book is available from the British Library

ISBN: HB: 978-1-5266-0021-9; TPB: 978-1-5266-0022-6; eBook: 978-1-5266-0019-6

2 4 6 8 10 9 7 5 3 1

Typeset by Newgen KnowledgeWorks Pvt. Ltd., Chennai, India
Printed and bound in Great Britain by CPI Group (UK) Ltd, Croydon CR0 4YY



To find out more about our authors and books visit www.bloomsbury.com
and sign up for our newsletters

*To the hundred billion people before us,
who fashioned our civilisation;
To the seven billion now alive,
whose actions may determine its fate;
To the trillions to come,
whose existence lies in the balance.*

PART ONE

THE STAKES

INTRODUCTION

If all goes well, human history is just beginning. Humanity is about two hundred thousand years old. But the Earth will remain habitable for hundreds of millions more—enough time for millions of future generations; enough to end disease, poverty and injustice forever; enough to create heights of flourishing unimaginable today. And if we could learn to reach out further into the cosmos, we could have more time yet: trillions of years, to explore billions of worlds. Such a lifespan places present-day humanity in its earliest infancy. A vast and extraordinary adulthood awaits.

Our view of this potential is easily obscured. The latest scandal draws our outrage; the latest tragedy, our sympathy. Time and space shrink. We forget the scale of the story in which we take part. But there are moments when we remember—when our vision shifts, and our priorities realign. We see a species precariously close to self-destruction, with a future of immense promise hanging in the balance. And which way that balance tips becomes our most urgent public concern.

This book argues that safeguarding humanity's future is the defining challenge of our time. For we stand at a crucial moment in the history of our species. Fuelled by technological progress, our power has grown so great that for the first time in humanity's long history, we have the capacity to destroy ourselves—severing our entire future and everything we could become.

Yet humanity's wisdom has grown only falteringly, if at all, and lags dangerously behind. Humanity lacks the maturity, coordination and foresight necessary to avoid making mistakes from

which we could never recover. As the gap between our power and our wisdom grows, our future is subject to an ever-increasing level of risk. This situation is unsustainable. So over the next few centuries, humanity will be tested: it will either act decisively to protect itself and its longterm potential, or, in all likelihood, this will be lost forever.

To survive these challenges and secure our future, we must act now: managing the risks of today, averting those of tomorrow, and becoming the kind of society that will never pose such risks to itself again.

It is only in the last century that humanity's power to threaten its entire future became apparent. One of the most harrowing episodes has just recently come to light. On Saturday 27 October 1962 a single officer on a Soviet submarine almost started a nuclear war. His name was Valentin Savitsky. He was captain of the submarine B-59—one of four submarines the Soviet Union had sent to support its military operations in Cuba. Each was armed with a secret weapon: a nuclear torpedo with explosive power comparable to the Hiroshima bomb.

It was the height of the Cuban Missile Crisis. Two weeks earlier, US aerial reconnaissance had produced photographic evidence that the Soviet Union was installing nuclear missiles in Cuba, from which they could strike directly at the mainland United States. In response, the US blockaded the seas around Cuba, drew up plans for an invasion and brought its nuclear forces to the unprecedented alert level of DEFCON 2 ('Next step to nuclear war').

On that Saturday, one of the blockading US warships detected Savitsky's submarine and attempted to force it to the surface by dropping low-explosive depth charges as warning shots. The submarine had been hiding deep underwater for days. It was out of radio contact, so the crew did not know whether war had already broken out. Conditions on board were extremely bad. It was built for the Arctic and its ventilator had broken in the tropical water. The heat inside was unbearable, ranging from 45 °C near the torpedo tubes to 60 °C in the engine room. Carbon dioxide had built

up to dangerous concentrations, and crew members had begun to fall unconscious. Depth charges were exploding right next to the hull. One of the crew later recalled: ‘It felt like you were sitting in a metal barrel, which somebody is constantly blasting with a sledgehammer.’

Increasingly desperate, Captain Savitsky ordered his crew to prepare their secret weapon:

Maybe the war has already started up there, while we are doing somersaults here. We’re going to blast them now! We will die, but we will sink them all—we will not disgrace our Navy!¹

Firing the nuclear weapon required the agreement of the submarine’s political officer, who held the other half of the firing key. Despite the lack of authorisation by Moscow, the political officer gave his consent.

On any of the other three submarines, this would have sufficed to launch their nuclear weapon. But by the purest luck, submarine B-59 carried the commander of the entire flotilla, Captain Vasili Arkhipov, and so required his additional consent. Arkhipov refused to grant it. Instead, he talked Captain Savitsky down from his rage and convinced him to give up: to surface amidst the US warships and await further orders from Moscow.²

We do not know precisely what would have happened if Arkhipov had granted his consent—or had he simply been stationed on any of the other three submarines. Perhaps Savitsky would not have followed through on his command. What is clear is that we came precariously close to a nuclear strike on the blockading fleet—a strike which would most likely have resulted in nuclear retaliation, then escalation to a full-scale nuclear war (the only kind the US had plans for). Years later, Robert McNamara, Secretary of Defense during the crisis, came to the same conclusion:

No one should believe that had U.S. troops been attacked by nuclear warheads, the U.S. would have refrained from responding with nuclear warheads. Where would it have ended? In utter disaster.³

Ever since the advent of nuclear weapons, humans have been making choices with such stakes. Ours is a world of flawed decision-makers, working with strikingly incomplete information, directing technologies which threaten the entire future of the species. We were lucky, that Saturday in 1962, and have so far avoided catastrophe. But our destructive capabilities continue to grow, and we cannot rely on luck forever.

We need to take decisive steps to end this period of escalating risk and safeguard our future. Fortunately, it is in our power to do so. The greatest risks are caused by human action, and they can be addressed by human action. Whether humanity survives this era is thus a choice humanity will make. But it is not an easy one. It all depends on how quickly we can come to understand and accept the fresh responsibilities that come with our unprecedented power.

This is a book about *existential risks*—risks that threaten the destruction of humanity’s longterm potential. Extinction is the most obvious way humanity’s entire potential could be destroyed, but there are others. If civilisation across the globe were to suffer a truly unrecoverable collapse, that too would destroy our longterm potential. And we shall see that there are dystopian possibilities as well: ways we might get locked into a failed world with no way back.

While this set of risks is diverse, it is also exclusive. So I will have to set aside many important risks that fall short of this bar: our topic is not new dark ages for humanity or the natural world (terrible though they would be), but the permanent destruction of humanity’s potential.

Existential risks present new kinds of challenges. They require us to coordinate globally and intergenerationally, in ways that go beyond what we have achieved so far. And they require foresight rather than trial and error. Since they allow no second chances, we need to build institutions to ensure that across our entire future we never once fall victim to such a catastrophe.

To do justice to this topic, we will have to cover a great deal of ground. Understanding the risks requires delving into physics, biology, earth science and computer science; situating this in the larger story of humanity requires history and anthropology; discerning just how much is at stake requires moral philosophy and economics; and finding solutions requires international relations and political science. Doing this properly requires deep engagement with each of these disciplines, not just cherry-picking expert quotes or studies that support one's preconceptions. This would be an impossible task for any individual, so I am extremely grateful for the extensive advice and scrutiny of dozens of the world's leading researchers from across these fields.⁴

This book is ambitious in its aims. Through careful analysis of the potential of humanity and the risks we face, it makes the case that we live during the most important era of human history. Major risks to our entire future are a new problem, and our thinking has not caught up. So *The Precipice* presents a new ethical perspective: a major reorientation in the way we see the world, and our role in it. In doing so, the book aspires to start closing the gap between our wisdom and power, allowing humanity a clear view of what is at stake, so that we will make the choices necessary to safeguard our future.

I have not always been focused on protecting our longterm future, coming to the topic only reluctantly. I am a philosopher, at Oxford University, specialising in ethics. My earlier work was rooted in the more tangible concerns of global health and global poverty—in how we could best help the worst off. When coming to grips with these issues I felt the need to take my work in ethics beyond the ivory tower. I began advising the World Health Organization, World Bank and UK government on the ethics of global health. And finding that my own money could do hundreds of times as much good for those in poverty as it could do for me, I made a lifelong pledge to donate at least a tenth of all I earn to help them.⁵ I founded a society, *Giving What We Can*, for those who wanted to join me, and was heartened to see thousands of

people come together to pledge more than £1 billion over our lifetimes to the most effective charities we know of, working on the most important causes. Together, we've already been able to transform the lives of tens of thousands of people.⁶ And because there are many other ways beyond our donations in which we can help fashion a better world, I helped start a wider movement, known as *effective altruism*, in which people aspire to use evidence and reason to do as much good as possible.

Since there is so much work to be done to fix the needless suffering in our present, I was slow to turn to the future. It was so much less visceral; so much more abstract. Could it really be as urgent a problem as suffering now? As I reflected on the evidence and ideas that would culminate in this book, I came to realise that the risks to humanity's future are just as real and just as urgent—yet even more neglected. And that the people of the future may be even more powerless to protect themselves from the risks we impose than the dispossessed of our own time.

Addressing these risks has now become the central focus of my work: both researching the challenges we face, and advising groups such as the UK Prime Minister's Office, the World Economic Forum and DeepMind on how they can best address these challenges. Over time, I've seen a growing recognition of these risks, and of the need for concerted action.

To allow this book to reach a diverse readership, I've been ruthless in stripping out the jargon, needless technical detail and defensive qualifications typical of academic writing (my own included). Readers hungry for further technical detail or qualifications can delve into the many endnotes and appendices, written with them in mind.⁷

I have tried especially hard to examine the evidence and arguments carefully and even-handedly, making sure to present the key points even if they cut against my narrative. For it is of the utmost importance to get to the truth of these matters—humanity's attention is scarce and precious, and must not be wasted on flawed narratives or ideas.⁸

Each chapter of *The Precipice* illuminates the central questions from a different angle. Part One (The Stakes) starts with a bird's-eye view of our unique moment in history, then examines why it warrants such urgent moral concern. Part Two (The Risks) delves into the science of the risks facing humanity, both from nature and from ourselves, showing that while some have been overstated, there is real risk and it is growing. So Part Three (The Path Forward) develops tools for understanding how these risks compare and combine, and new strategies for addressing them. I close with a vision of our future: of what we could achieve were we to succeed.

This book is not just a familiar story of the perils of climate change or nuclear war. These risks that first awoke us to the possibilities of destroying ourselves are just the beginning. There are emerging risks, such as those arising from biotechnology and advanced artificial intelligence, that may pose much greater risk to humanity in the coming century.

Finally, this is not a pessimistic book. It does not present an inevitable arc of history culminating in our destruction. It is not a morality tale about our technological hubris and resulting fall. Far from it. The central claim is that there are real risks to our future, but that our choices can still make all the difference. I believe we are up to the task: that through our choices we can pull back from the precipice and, in time, create a future of astonishing value—with a richness of which we can barely dream, made possible by innovations we are yet to conceive. Indeed, my deep optimism about humanity's future is core to my motivation in writing this book. Our potential is vast. We have so much to protect.

1

STANDING AT THE PRECIPICE

It might be a familiar progression, transpiring on many worlds—a planet, newly formed, placidly revolves around its star; life slowly forms; a kaleidoscopic procession of creatures evolves; intelligence emerges which, at least up to a point, confers enormous survival value; and then technology is invented. It dawns on them that there are such things as laws of Nature, that these laws can be revealed by experiment, and that knowledge of these laws can be made both to save and to take lives, both on unprecedented scales. Science, they recognize, grants immense powers. In a flash, they create world-altering contrivances. Some planetary civilizations see their way through, place limits on what may and what must not be done, and safely pass through the time of perils. Others, not so lucky or so prudent, perish.

—Carl Sagan¹

We live at a time uniquely important to humanity's future. To see why, we need to take a step back and view the human story as a whole: how we got to this point and where we might be going next.

Our main focus will be humanity's ever-increasing power—power to improve our condition and power to inflict harm. We shall see how the major transitions in human history have enhanced our power, and enabled us to make extraordinary progress. If we can avoid catastrophe we can cautiously expect this progress to continue: the future of a responsible humanity is

extraordinarily bright. But this increasing power has also brought on a new transition, at least as significant as any in our past, the transition to our time of perils.

HOW WE GOT HERE

Very little of humanity's story has been told; because very little *can* be told. Our species, *Homo sapiens*, arose on the savannahs of Africa 200,000 years ago.² For an almost unimaginable time we have had great loves and friendships, suffered hardships and griefs, explored, created, and wondered about our place in the universe. Yet when we think of humanity's great achievements across time, we think almost exclusively of deeds recorded on clay, papyrus or paper—records that extend back only about 5,000 years. We rarely think of the first person to set foot in the strange new world of Australia some 70,000 years ago; of the first to name and study the plants and animals of each place we reached; of the stories, songs and poems of humanity in its youth.³ But these accomplishments were real, and extraordinary.

We know that even before agriculture or civilisation, humanity was a fresh force in the world. Using the simple, yet revolutionary, technologies of seafaring, clothing and fire, we travelled further than any mammal before us. We adapted to a wider range of environments, and spread across the globe.⁴

What made humanity exceptional, even at this nascent stage? We were not the biggest, the strongest or the hardest. What set us apart was not physical, but mental—our intelligence, creativity and language.⁶

Yet even with these unique mental abilities, a single human alone in the wilderness would be nothing exceptional. He or she might be able to survive—intelligence making up for physical prowess—but would hardly dominate. In ecological terms, it is not a *human* that is remarkable, but *humanity*.

Each human's ability to cooperate with the dozens of other people in their band was unique among large animals. It allowed us to form something greater than ourselves. As our language

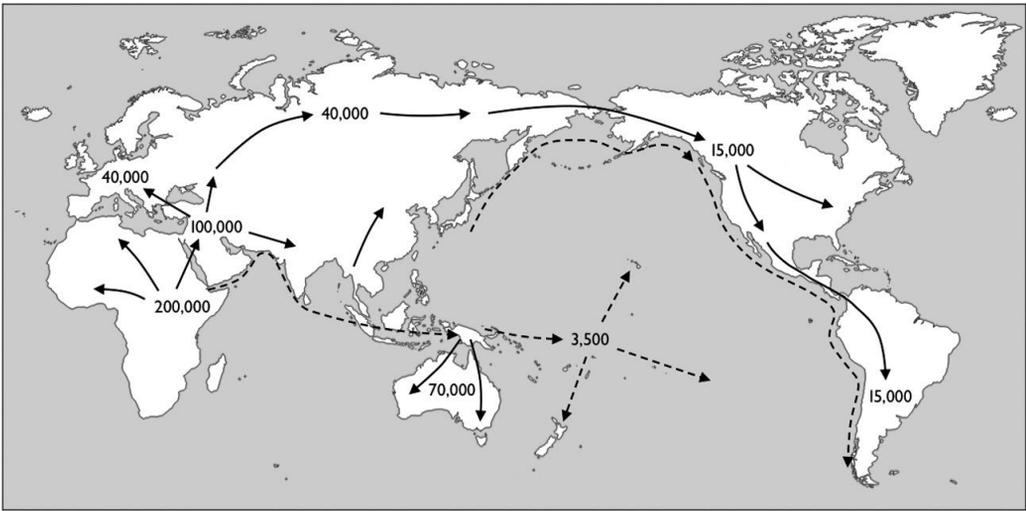


FIGURE 1.1 How we settled the world. The arrows show our current understanding of the land and sea routes taken by our ancestors, and how many years ago they reached each area.⁵

grew in expressiveness and abstraction, we were able to make the most of such groupings: pooling together our knowledge, our ideas and our plans.

Crucially, we were able to cooperate across *time* as well as space. If each generation had to learn everything anew, then even a crude iron shovel would have been forever beyond our technological reach. But we learned from our ancestors, added minor innovations of our own, and passed this all down to our children. Instead of dozens of humans in cooperation, we had tens of thousands, cooperating across the generations, preserving and improving ideas through deep time. Little by little, our knowledge and our culture grew.⁷

At several points in the long history of humanity there has been a great transition: a change in human affairs that accelerated our accumulation of power and shaped everything that would follow. I will focus on three.⁸

The first was the Agricultural Revolution.⁹ Around 10,000 years ago the people of the Fertile Crescent, in the Middle East, began planting wild wheat, barley, lentils and peas to supplement their foraging. By preferentially replanting the

seeds from the best plants, they harnessed the power of evolution, creating new domesticated varieties with larger seeds and better yields. This worked with animals too, giving humans easier access to meat and hides, along with milk, wool and manure. And the physical power of draft animals to help plough the fields or transport the harvest was the biggest addition to humanity's power since fire.¹⁰

While the Fertile Crescent is often called 'the cradle of civilisation', in truth civilisation had many cradles. Entirely independent agricultural revolutions occurred across the world in places where the climate and local species were suitable: in east Asia; sub-Saharan Africa; New Guinea; South, Central and North America; and perhaps elsewhere too.¹¹ The new practices fanned out from each of these cradles, changing the way of life for many from foraging to farming.

This had dramatic effects on the scale of human cooperation. Agriculture reduced the amount of land needed to support each person by a factor of a hundred, allowing large permanent settlements to develop, which began to unite together into states.¹² Where the largest foraging communities involved perhaps hundreds of people, some of the first cities had tens of thousands of inhabitants. At its height, the Sumerian civilisation contained around a million people.¹³ And 2,000 years ago, the Han dynasty of China reached sixty million people—about a *hundred thousand* times as many as were ever united in our forager past, and about ten times the entire global forager population at its peak.¹⁴

As more and more people were able to share their insights and discoveries, there were rapid developments in technology, institutions and culture. And the increasing numbers of people trading with one another made it possible for them to specialise in these areas—to devote a lifetime to governance, trade or the arts—allowing us to develop these ideas much more deeply.

Over the first 6,000 years of agriculture, we achieved world-changing breakthroughs including writing, mathematics, law and the wheel.¹⁵ Of these, writing was especially important



FIGURE 1.2 The cradles of civilisation. The places around the world where agriculture was independently developed, marked with how many years ago this occurred.

for strengthening our ability to cooperate across time and space: increasing the bandwidth between generations, the reliability of the information, and the distance over which ideas could be shared.

The next great transition was the Scientific Revolution.¹⁶ Early forms of science had been practised since ancient times, and the seeds of empiricism can be found in the work of medieval scholars in the Islamic world and Europe.¹⁷ But it was only about 400 years ago that humanity developed the scientific method and saw scientific progress take off.¹⁸ This helped replace a reliance on received authorities with careful observation of the natural world, seeking simple and testable explanations for what we saw. The ability to test and discard bad explanations helped us break free from dogma, and allowed for the first time the systematic creation of knowledge about the workings of nature.

Some of our new-found knowledge could be harnessed to improve the world around us. So the accelerated accumulation of knowledge brought with it an acceleration of technological innovation, giving humanity increasing power over the natural world. The rapid pace allowed people to see transformative effects of these improvements within their own lifetimes. This gave rise to

the modern idea of *progress*. Where the world had previously been dominated by narratives of decline and fall or of a recurring cycle, there was increasing interest in a new narrative: a grand project of working together to build a better future.

Soon, humanity underwent a third great transition: the Industrial Revolution. This was made possible by the discovery of immense reserves of energy in the form of coal and other fossil fuels. These are formed from the compressed remains of organisms that lived in aeons past, allowing us access to a portion of the sunlight that shone upon the Earth over millions of years.¹⁹ We had already begun to drive simple machines with the renewable energy from the wind, rivers and forests; fossil fuels allowed access to vastly more energy, and in a much more concentrated and convenient form.

But energy is nothing without a way of converting it to useful work, to achieve our desired changes in the world. The steam engine allowed the stored chemical energy of coal to be turned into mechanical energy.²⁰ This mechanical energy was then used to drive machines that performed massive amounts of labour for us, allowing raw materials to be transformed into finished products much more quickly and cheaply than before. And via the railroad, this wealth could be distributed and traded across long distances.

Productivity and prosperity began to accelerate, and a rapid sequence of innovations ramped up the efficiency, scale and variety of automation, giving rise to the modern era of sustained economic growth.²¹

The effects of these transitions have not always been positive. Life in the centuries following the Agricultural Revolution generally involved more work, reduced nutrition and increased disease.²² Science gave us weapons of destruction that haunt us to this day. And the Industrial Revolution was among the most destabilising periods in human history. The unequal distribution of gains in prosperity and the exploitative labour practices led to the revolutionary upheavals of the early twentieth century.²³

Inequality between countries increased dramatically (a trend that has only begun to reverse in the last two decades).²⁴ Harnessing the energy stored in fossil fuels has released greenhouse gases, while industry fuelled by this energy has endangered species, damaged ecosystems and polluted our environment.

Yet despite these real problems, on average human life today is substantially better than at any previous time. The most striking change may be in breaking free from poverty. Until 200 years ago—the last thousandth of our history²⁵—increases in humanity's power and prosperity came hand in hand with increases in the human population. Income *per person* stayed almost unchanged: a little above subsistence in times of plenty; a little below in times of need.²⁶ The Industrial Revolution broke this rule, allowing income to grow faster than population and ushering in an unprecedented rise in prosperity that continues to this day.

We often think of economic growth from the perspective of a society that is already affluent, where it is not immediately clear if further growth even improves our lives. But the most remarkable effects of economic growth have been for the poorest people. In today's world, one out of ten people are so poor that they live on less than two dollars per day—a widely used threshold for 'extreme poverty'. That so many have so little is among the greatest problems of our time, and has been a major focus of my life. It is shocking then to look further back and see that prior to the Industrial Revolution 19 out of 20 people lived on less than two dollars a day (even adjusting for inflation and purchasing power). Until the Industrial Revolution, any prosperity was confined to a tiny elite with extreme poverty the norm. But over the last two centuries more and more people have broken free from extreme poverty, and are now doing so more quickly than at any earlier time.²⁷ Two dollars a day is far from prosperity, and these statistics can be of little comfort to those who are still in the grip of poverty, but the trends towards improvement are clear.

And it is not only in terms of material conditions that life has improved. Consider education and health. Universal schooling

has produced dramatic improvements in education. Before the Industrial Revolution, just one in ten of the world's people could read and write; now more than eight in ten can do so.²⁸ For the 10,000 years since the Agricultural Revolution, life expectancy had hovered between 20 and 30 years. It has now more than doubled, to 72 years.²⁹ And like literacy, these gains have been felt across the world. In 1800 the highest life expectancy of any country was a mere 43 years, in Iceland. Now every single country has a life expectancy above 50.³⁰ The industrial period has seen all of humanity become more prosperous, educated and long-lived than ever before. But we should not succumb to complacency in the face of this astonishing progress. That we have achieved so much, and so quickly, should inspire us to address the suffering and injustices that remain.

We have also seen substantial improvements in our moral thinking.³² One of the clearest trends is towards the gradual expansion of the moral community, with the recognition of the rights of women, children, the poor, foreigners and ethnic or religious minorities. We have also seen a marked shift away from violence as a morally acceptable part of society.³³ And in the last sixty years we have added the environment and the welfare of animals to our standard picture of morality. These social changes did not come naturally with prosperity. They were secured by reformers and activists, motivated by the belief that we can—and must—improve. We still have far to go before we are living up to these new ideals, and our progress can be painfully slow, but looking back even just one or two centuries shows how far we have come.

Of course, there have been many setbacks and exceptions. The path has been tumultuous, things have often become better in some ways while worse in others, and there is certainly a danger of choosing selectively from history to create a simple narrative of improvement from a barbarous past to a glorious present. Yet at the largest scales of human history, where we see not the rise and fall of each empire, but the changing face of human civilisation across the entire globe, the trends towards progress are clear.³⁴

It can be hard to believe such trends, when it so often feels like everything is collapsing around us. In part this scepticism comes from our everyday experience of our own lives or communities over a timespan of years—a scale where downs are almost as likely as ups. It might also come from our tendency to focus more on bad news than good and on threats rather than opportunities: heuristics that are useful for directing our actions, but which misfire when attempting to objectively assess the balance of bad and good.³⁵ When we try to overcome these distortions, looking for global indicators of the quality of our lives that are as objective as possible, it is very difficult to avoid seeing significant improvement from century to century.

And these trends should not surprise us. Every day we are the beneficiaries of uncountable innovations made by people over hundreds of thousands of years. Innovations in technology, mathematics, language, institutions, culture, art; the ideas of the hundred billion people who came before us, and shaped almost every facet of the modern world.³⁶ This is a stunning inheritance. No wonder, then, that our lives are better for it.

We cannot be sure these trends towards progress will continue. But given their tenacity, the burden would appear to be on the pessimist to explain why *now* is the point it will fail. This is especially true when people have been predicting such failure for so long and with such a poor track record. Thomas Macaulay made this point well:

We cannot absolutely prove that those are in error who tell us that society has reached a turning point, that we have seen our best days. But so said all before us, and with just as much apparent reason . . . On what principle is it that, when we see nothing but improvement behind us, we are to expect nothing but deterioration before us?³⁷

And he wrote those words in 1830, before an additional 190 years of progress and failed predictions of the end of progress. During those years, lifespan doubled, literacy soared and eight in ten people escaped extreme poverty. What might the coming years bring?

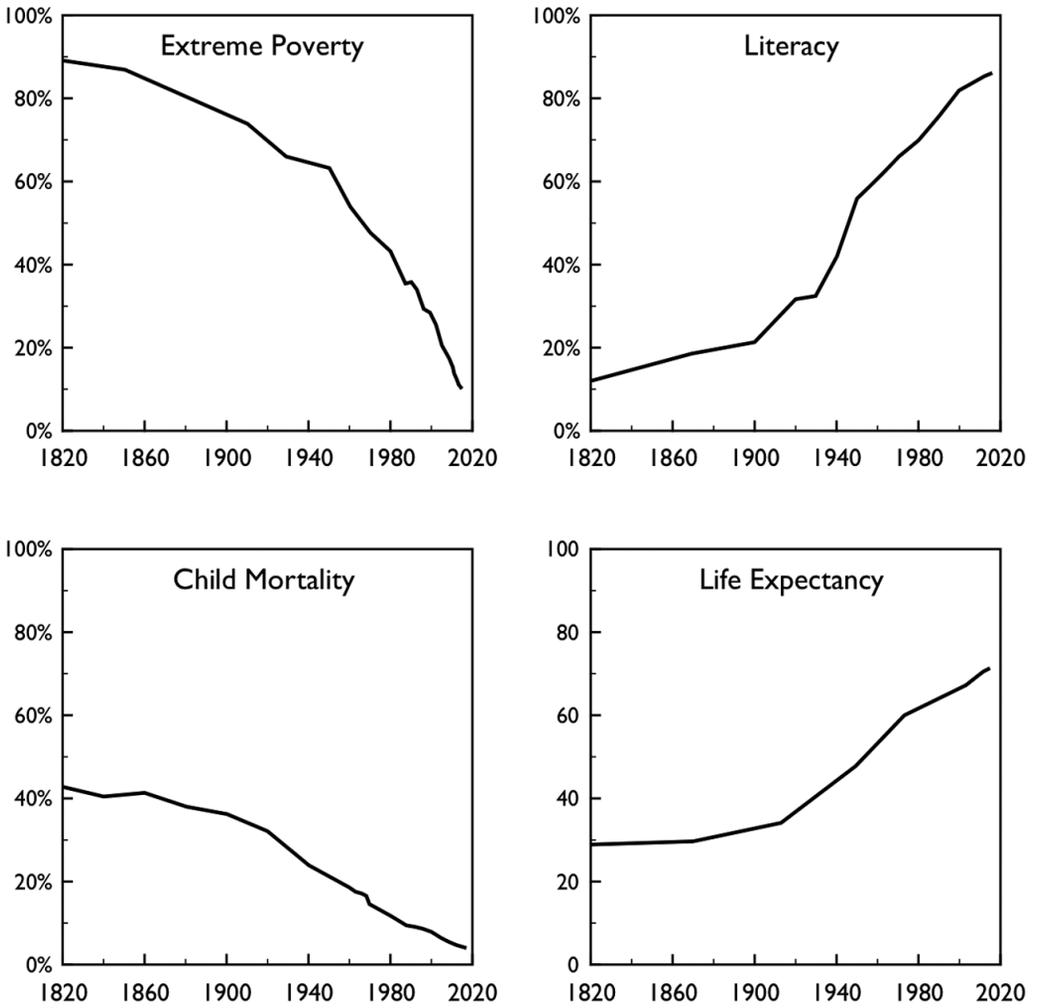


FIGURE 1.3 The striking improvements in extreme poverty, literacy, child mortality and life expectancy over the last 200 years.³¹

WHERE WE MIGHT GO

On the timescale of an individual human life, our 200,000-year history seems almost incomprehensibly long. But on a geological timescale it is short, and vanishingly so on the timescale of the universe as a whole. Our cosmos has a 14-billion-year history, and even that is short on the grandest scales. Trillions of years lie ahead of us. The future is immense.

How much of this future might we live to see? The fossil record provides some useful guidance. Mammalian species typically survive for around one million years before they go extinct; our

close relative, *Homo erectus*, survived for almost two million.³⁸ If we think of one million years in terms of a single, eighty-year life, then today humanity would be in its adolescence—sixteen years old; just coming into our power; just old enough to get ourselves in serious trouble.³⁹

Obviously, though, humanity is not a typical species. For one thing, we have recently acquired a unique power to destroy ourselves—power that will be the focus of much of this book. But we also have unique power to protect ourselves from external destruction, and thus the potential to outlive our related species.

How long *could* we survive on Earth? Our planet will remain habitable for roughly a billion years.⁴⁰ That's enough time for trillions of human lives; time to watch mountain ranges rise, continents collide, orbits realign; and time, as well, to heal our society and our planet of the wounds we have caused in our immaturity.

And we might have more time yet. As one of the pioneers of rocketry put it, 'Earth is the cradle of humanity, but one cannot live in a cradle forever.'⁴¹ We do not know, yet, how to reach other stars and settle their planets, but we know of no fundamental obstacles. The main impediment appears to be the time necessary to learn how. This makes me optimistic. After all, the first heavier-than-air flight was in 1903 and just sixty-eight years later we had launched a spacecraft that left our Solar System and will reach the stars. Our species learns quickly, especially in recent times, and a billion years is a long education. I think we will need far less.

If we can reach other stars, then the whole galaxy opens up to us. The Milky Way alone contains more than 100 billion stars, and some of these will last for trillions of years, greatly extending our potential lifespan. Then there are billions of other galaxies beyond our own. If we reach a future of such a scale, we might have a truly staggering number of descendants, with the time, resources, wisdom and experience to create a diversity of wonders unimaginable to us today.

While humanity has made progress towards greater prosperity, health, education and moral inclusiveness, there is so much further we could go. Our present world remains marred by malaria and HIV; depression and dementia; racism and sexism; torture and oppression. But with enough time, we can end these horrors—building a society that is truly just and humane.

And a world without agony and injustice is just a lower bound on how good life could be. Neither the sciences nor the humanities have yet found any upper bound. We get some hint at what is possible during life's best moments: glimpses of raw joy, luminous beauty, soaring love. Moments when we are truly awake. These moments, however brief, point to possible heights of flourishing far beyond the status quo, and far beyond our current comprehension.

Our descendants could have aeons to explore these heights, with new means of exploration. And it's not just wellbeing. Whatever you value—beauty, understanding, culture, consciousness, freedom, adventure, discovery, art—our descendants would be able to take these so much further, perhaps even discovering entirely new categories of value, completely unknown to us. Music we lack the ears to hear.

THE PRECIPICE

But this future is at risk. For we have recently undergone another transition in our power to transform the world—one at least as significant as the Agricultural, Scientific and Industrial Revolutions that preceded it.

With the detonation of the first atomic bomb, a new age of humanity began.⁴² At that moment, our rapidly accelerating technological power finally reached the threshold where we might be able to destroy ourselves. The first point where the threat to humanity from within exceeded the threats from the natural world. A point where the entire future of humanity hangs in the balance. Where every advance our ancestors have made could be squandered, and every advance our descendants may achieve

could be denied. The greater part of the book of human history left unwritten; the narrative broken off; blank pages.

Nuclear weapons were a discontinuous change in human power. At Hiroshima, a single bomb did the damage of thousands. And six years later, a single thermonuclear bomb held more energy than every explosive used in the entire course of the Second World War.⁴³

It became clear that a war with such weapons would change the Earth in ways that were unprecedented in human history. World leaders, atomic scientists and public intellectuals began to take seriously the possibility that a nuclear war would spell the end of humanity: either through extinction or a permanent collapse of civilisation.⁴⁴ Early concern centred on radioactive fallout and damage to the ozone layer, but in the 1980s the focus shifted to a scenario known as nuclear winter, in which nuclear firestorms loft smoke from burning cities into the upper atmosphere.⁴⁵ High above the clouds, the smoke cannot be rained out and would persist for years, blackening the sky, chilling the Earth and causing massive crop failure. This was a mechanism by which nuclear war could result in extreme famine, not just in the combatant countries, but in every country around the world. Millions of direct deaths from the explosions could be followed by billions of deaths from starvation, and—potentially—by the end of humanity itself.

How close have we come to such a war? With so much to lose, nuclear war is in no one's interest. So we might expect these obvious dangers to create a certain kind of safety—where world leaders inevitably back down before the brink. But as more and more behind-the-scenes evidence from the Cold War has become public, it has become increasingly clear that we have only barely avoided full-scale nuclear war.

We saw how the intervention of a single person, Captain Vasili Arkhipov, may have prevented an all-out nuclear war at the height of the Cuban Missile Crisis. But even more shocking is

just how many times in those few days we came close to disaster, only to be pulled back by the decisions of a few individuals.

The principal events of the crisis took place over a single week. On Monday 22 October 1962, President John F. Kennedy gave a television address, informing his nation that the Soviets had begun installing strategic nuclear missiles in Cuba—directly threatening the United States. He warned that any use of these nuclear weapons would be met by a full-scale nuclear retaliation on the Soviet Union. His advisors drew up plans for both air strikes on the 48 missiles they had discovered and a full invasion of Cuba. US forces were brought to DEFCON 3, to prepare for a possible nuclear war.⁴⁶

On Wednesday 24 October the US launched a naval blockade to prevent the delivery of further missiles to Cuba, and took its nuclear forces to the unprecedented level of DEFCON 2. Nuclear missiles were readied for launch and nuclear bombers took to the skies, ready to begin an all-out nuclear attack on the Soviet Union. The crisis reached its peak on Saturday when the Soviets shot down a U-2 reconnaissance plane with a surface-to-air missile, killing its pilot.

Then on Sunday morning it was all over. The Soviets backed down, unexpectedly announcing that they were removing all nuclear missiles from Cuba. But it could very easily have ended differently.

There has been substantial debate about exactly how close the crisis came to nuclear war. But over the decades, as more details have been revealed, the picture has become increasingly serious. Kennedy and Khrushchev went to great lengths to resist hawkish politicians and generals and to stay clear of the brink.⁴⁷ But there was a real possibility that, like the First World War, a war might begin without any side wanting it. As the week wore on, events on the ground spiralled beyond their control and they only barely kept the crisis from escalating. The US came extremely close to attacking Cuba, this had a much higher chance of causing nuclear retaliation than anyone guessed, and this in turn had a high chance of escalating to full-scale nuclear war.

Twice, during the crisis, the US nearly launched an attack on Cuba. At the height of the tensions, Kennedy had agreed that if a U-2 were shot down, the US would immediately strike Cuba, with no need to reconvene the war council. Then, on Saturday, a U-2 was indeed shot down. But Kennedy changed his mind and called off the counter-attack. Instead, he issued a secret ultimatum, informing the Soviets that if they did not commit to removing the missiles within twenty-four hours, or if another plane was shot down, the US would immediately launch air strikes and, almost surely, a full invasion.

This too almost triggered an attack. For the Americans did not know the extent to which Khrushchev was unable to control his forces in Cuba. Indeed, the U-2 had been shot down by a Soviet general acting against explicit orders from Khrushchev. And Khrushchev had even less control over the Cuban forces, who had already hit a low-flying reconnaissance plane with anti-aircraft fire and were eager to take one down. Knowing that he could not stop his own side from downing another plane, thereby triggering a US attack, Khrushchev raced to issue a statement ending the crisis before morning reconnaissance flights resumed.

What would have happened if the US *had* attacked? American leaders assumed that a purely conventional (non-nuclear) attack on Cuba could only be met with a purely conventional response. It was out of the question, they thought, that the Soviets would respond with nuclear attacks on the mainland United States. But they were missing another crucial fact. The missiles the US had discovered in Cuba were only a fraction of those the Soviets had delivered. There were 158 nuclear warheads. And more than 90 of these were tactical nuclear weapons, there for the express purpose of nuclear first use: to destroy a US invasion fleet before it could land.⁴⁸

What's more, Castro was eager to use them. Indeed, he directly asked Khrushchev to fire the nuclear weapons if the Americans tried to invade, even though he knew this would lead to the annihilation of his own country: 'What would have happened to Cuba? It would have been totally destroyed.'⁴⁹ And

Khrushchev, in another unprecedented move, had relinquished central control of the tactical nuclear weapons, delegating the codes and decision to fire to the local Soviet commander. After hearing Kennedy's television address, Khrushchev issued new orders that the weapons were not to be used without his explicit permission, but he came to fear these would be disobeyed in the heat of conflict, as his order not to fire on US spy planes had been.

So unbeknownst to the US military leadership, a conventional attack on Cuba was likely to be met with a nuclear strike on American forces. And such a strike was extremely likely to be met by a further nuclear response from the US. This nuclear response was highly likely to go beyond Cuba, and to precipitate a full-scale nuclear war with the Soviets. In his television address on the Monday, Kennedy had explicitly promised that 'It shall be the policy of this Nation to regard any nuclear missile launched from Cuba against any nation in the Western Hemisphere as an attack by the Soviet Union on the United States, requiring a full retaliatory response upon the Soviet Union.'⁵⁰

It is extremely difficult to estimate the chance that the crisis would have escalated to nuclear war.⁵¹ Shortly after, Kennedy told a close advisor that he thought the probability of it ending in nuclear war with the USSR was 'somewhere between one out of three, and even'.⁵² And it has just been revealed that the day after the crisis ended, Paul Nitze (an advisor to Kennedy's war council) estimated the chance at 10 percent, and thought that everyone else in the council would have put it even higher.⁵³ Moreover, none of these people knew about the tactical nuclear weapons in Cuba, Khrushchev's lack of control of his troops or the events on submarine B-59.

While I'm reluctant to question those whose very decisions could have started the war, my own view is that they were somewhat too pessimistic, given what they knew at the time. However, when we include the subsequent revelations about what was

really happening in Cuba my estimates would roughly match theirs. I'd put the chance of the crisis escalating to a nuclear war with the Soviets at something between 10 and 50 percent.⁵⁴

When writing about such close calls, there is a tendency to equate this chance to that of the end of civilisation or the end of humanity itself. But that would be a large and needless exaggeration. For we need to combine this chance of nuclear war with the chance that such a war would spell the end of humanity or human civilisation, which is far from certain. Yet even making such allowances the Cuban Missile Crisis would remain one of the pivotal moments in 200,000 years of human history: perhaps the closest we have ever come to losing it all.

Even now, with the Cold War just a memory, nuclear weapons still pose a threat to humanity. At the time of writing, the highest chance of a nuclear conflict probably involves North Korea. But not all nuclear wars are equal. North Korea has less than 1 percent as many warheads as Russia or the USA, and they are substantially smaller. A nuclear war with North Korea would be a terrible disaster, but it currently poses little threat to humanity's longterm potential.⁵⁵

Instead, most of the existential risk from nuclear weapons today probably still comes from the enormous American and Russian arsenals. The development of ICBMs (intercontinental ballistic missiles) allowed each side to destroy most of the other's missiles with just thirty minutes' warning, so they each moved many missiles to 'hair-trigger alert'—ready to launch in just ten minutes.⁵⁶ Such hair-trigger missiles are extremely vulnerable to accidental launch, or to deliberate launch during a false alarm. As we shall see in Chapter 4, there has been a chilling catalogue of false alarms continuing past the end of the Cold War. On a longer timescale there is also the risk of other nations creating their own enormous stockpiles, of innovations in military technologies undermining the logic of deterrence, and of shifts in the geopolitical landscape igniting another arms race between great powers.

Nuclear weapons are not the only threat to humanity. They have been our focus so far because they were the first major risk and have already threatened humanity. But there are others too.

The exponential rise in prosperity brought on by the Industrial Revolution came on the back of a rapid rise in carbon emissions. A minor side effect of industrialisation has eventually grown to become a global threat to health, the environment, international stability, and maybe even humanity itself.

Nuclear weapons and climate change have striking similarities and contrasts. They both threaten humanity through major shifts in the Earth's temperature, but in opposite directions. One burst in upon the scene as the product of an unpredictable scientific breakthrough; the other is the continuation of centuries-long scaling-up of old technologies. One poses a small risk of sudden and precipitous catastrophe; the other is a gradual, continuous process, with a delayed onset—where some level of catastrophe is assured and the major uncertainty lies in just how bad it will be. One involves a classified military technology controlled by a handful of powerful actors; the other involves the aggregation of small effects from the choices of everyone in the world.

As technology continues to advance, new threats appear on the horizon. These threats promise to be more like nuclear weapons than like climate change: resulting from sudden breakthroughs, precipitous catastrophes, and the actions of a small number of actors. There are two emerging technologies that especially concern me; they will be the focus of Chapter 5.

Ever since the Agricultural Revolution, we have induced genetic changes in the plants and animals around us to suit our ends. But the discovery of the genetic code and the creation of tools to read and write it have led to an explosion in our ability to refashion life to new purposes. Biotechnology will bring major improvements in medicine, agriculture and industry. But it will also bring risks to civilisation and to humanity itself: both from accidents during legitimate research and from engineered bioweapons.

We are also seeing rapid progress in the capabilities of AI systems, with the biggest improvements in the areas where AI has

traditionally been weakest, such as perception, learning and general intelligence. Experts find it likely that this will be the century that AI exceeds human ability not just in a narrow domain, but in general intelligence—the ability to overcome a diverse range of obstacles to achieve one’s goals. Humanity has risen to a position where we control the rest of the world precisely because of our unparalleled mental abilities. If we pass this mantle to our machines, it will be they who are in this unique position. This should give us cause to wonder why it would be humanity who will continue to call the shots. We need to learn how to align the goals of increasingly intelligent and autonomous machines with human interests, and we need to do so before those machines become more powerful than we are.

These threats to humanity, and how we address them, define our time. The advent of nuclear weapons posed a real risk of human extinction in the twentieth century. With the continued acceleration of technology, and without serious efforts to protect humanity, there is strong reason to believe the risk will be higher this century, and increasing with each century that technological progress continues. Because these anthropogenic risks outstrip all natural risks combined, they set the clock on how long humanity has left to pull back from the brink.

I am not claiming that extinction is the inevitable conclusion of scientific progress, or even the most likely outcome. What I am claiming is that there has been a robust trend towards increases in the power of humanity which has reached a point where we pose a serious risk to our own existence. How we react to this risk is up to us.

Nor am I arguing against technology. Technology has proved itself immensely valuable in improving the human condition. And technology is essential for humanity to achieve its longterm potential. Without it, we would be doomed by the accumulated risk of natural disasters such as asteroid impacts. Without it, we would never achieve the highest flourishing of which we are capable.

The problem is not so much an excess of technology as a lack of wisdom.⁵⁷ Carl Sagan put this especially well:

Many of the dangers we face indeed arise from science and technology—but, more fundamentally, because we have become powerful without becoming commensurately wise. The world-altering powers that technology has delivered into our hands now require a degree of consideration and foresight that has never before been asked of us.⁵⁸

This idea has even been advocated by a sitting US president:

the very spark that marks us as a species—our thoughts, our imagination, our language, our tool-making, our ability to set ourselves apart from nature and bend it to our will—those very things also give us the capacity for unmatched destruction . . . Technological progress without an equivalent progress in human institutions can doom us. The scientific revolution that led to the splitting of an atom requires a moral revolution as well.⁵⁹

We need to gain this wisdom; to have this moral revolution. Because we cannot come back from extinction, we cannot wait until a threat strikes before acting—we must be proactive. And because gaining wisdom or starting a moral revolution takes time, we need to start now.

I think that we are likely to make it through this period. Not because the challenges are small, but because we will rise to them. The very fact that these risks stem from human action shows us that human action can address them.⁶⁰ Defeatism would be both unwarranted and counterproductive—a self-fulfilling prophecy. Instead, we must address these challenges head-on with clear and rigorous thinking, guided by a positive vision of the longterm future we are trying to protect.

How big are these risks? One cannot expect precise numbers, as the risks are *complex* (so not amenable to simple mathematical analysis) and *unprecedented* (so cannot be approximated by a longterm frequency). Yet it is important to at least try

to give quantitative estimates. Qualitative statements such as ‘a grave risk of human extinction’ could be interpreted as meaning anything from 1 percent all the way to 99 percent.⁶¹ They add more confusion than clarity. So I will offer quantitative estimates, with the proviso that they can’t be precise and are open to revision.

During the twentieth century, my best guess is that we faced around a one in a hundred risk of human extinction or the unrecoverable collapse of civilisation. Given everything I know, I put the existential risk this century at around one in six: Russian roulette.⁶² (See table 6.1 on p. 167 for a breakdown of the risks.) If we do not get our act together, if we continue to let our growth in power outstrip that of wisdom, we should expect this risk to be even higher next century, and each successive century.

These are the greatest risks we have faced.⁶³ If I’m even roughly right about their scale, then we cannot survive many centuries with risk like this. It is an *unsustainable* level of risk.⁶⁴ Thus, one way or another, this period is unlikely to last more than a small number of centuries.⁶⁵ Either humanity takes control of its destiny and reduces the risk to a sustainable level, or we destroy ourselves.

Consider human history as a grand journey through the wilderness. There are wrong turns and times of hardship, but also times of sudden progress and heady views. In the middle of the twentieth century we came through a high mountain pass and found that the only route onward was a narrow path along the cliff-side: a crumbling ledge on the brink of a precipice. Looking down brings a deep sense of vertigo. If we fall, everything is lost. We do not know just how likely we are to fall, but it is the greatest risk to which we have ever been exposed.

This comparatively brief period is a unique challenge in the history of our species. Our response to it will define our story. Historians of the future will name this time, and schoolchildren will study it. But I think we need a name now. I call it the Precipice.

The Precipice gives our time immense meaning. In the grand course of history—if we make it that far—*this* is what our time

will be remembered for: for the highest levels of risk, and for humanity opening its eyes, coming into its maturity and guaranteeing its long and flourishing future. This is the meaning of our time.

I am not glorifying our generation, nor am I vilifying us. The point is that our actions have uniquely high stakes. Whether we are great or terrible will depend upon what we do with this opportunity. I hope we live to tell our children and grandchildren that we did not stand by, but used this chance to play the part that history gave us.

Safeguarding humanity through these dangers should be a central priority of our time. I am not saying that this is the only issue in the world, that people should drop everything else they hold dear to do this. But if you can see a way that you could play a role—if you have the skills, or if you are young and can shape your own path—then I think safeguarding humanity through these times is among the most noble purposes you could pursue.

THE PRECIPICE & ANTHROPOCENE

It has become increasingly clear that human activity is the dominant force shaping the environment. Scientists are concluding that humanity looms large not just in its own terms, but in the objective terms of biology, geology and climatology. If there are geologists in the distant future, they would identify the layer of rock corresponding to our time as a fundamental change from the layers before it.

Our present-day geologists are thus considering making this official—changing their classification of geological time to introduce a new epoch called the *Anthropocene*. Proposed beginnings for this epoch include the megafauna extinctions, the Agricultural Revolution, the crossing of the Atlantic, the Industrial Revolution, or early nuclear weapons tests.⁶⁶

Is this the same as the Precipice? How do they differ?

- The Anthropocene is the time of profound human effects on the environment, while the Precipice is the time where humanity is at high risk of destroying itself.
- The Anthropocene is a geological epoch, which typically last millions of years, while the Precipice is a time in human history (akin to the Enlightenment or the Industrial Revolution), which will likely end within a few centuries.
- They might both officially start with the first atomic test, but this would be for very different reasons. For the Anthropocene, it would be mainly for convenient dating, while for the Precipice it is because of the risk nuclear weapons pose to our survival.

Buy the book on Amazon:
<https://geni.us/the-precipice>

Other purchase options,
including audiobook:
<https://theprecipice.com/purchase>