

·1·

Introduction

Nick Bostrom and Milan M. Ćirković

1.1 Why?

The term ‘global catastrophic risk’ lacks a sharp definition. We use it to refer, loosely, to a risk that might have the potential to inflict serious damage to human well-being on a global scale. On this definition, an immensely diverse collection of events could constitute global catastrophes: potential candidates range from volcanic eruptions to pandemic infections, nuclear accidents to worldwide tyrannies, out-of-control scientific experiments to climatic changes, and cosmic hazards to economic collapse. With this in mind, one might well ask, what use is a book on global catastrophic risk? The risks under consideration seem to have little in common, so does ‘global catastrophic risk’ even make sense as a topic? Or is the book that you hold in your hands as ill-conceived and unfocused a project as a volume on ‘Gardening, Matrix Algebra, and the History of Byzantium’?

We are confident that a comprehensive treatment of global catastrophic risk will be at least somewhat more useful and coherent than the above-mentioned imaginary title. We also believe that studying this topic is highly important. Although the risks are of various kinds, they are tied together by many links and commonalities. For example, for many types of destructive events, much of the damage results from second-order impacts on social order; thus the risks of social disruption and collapse are not unrelated to the risks of events such as nuclear terrorism or pandemic disease. Or to take another example, apparently dissimilar events such as large asteroid impacts, volcanic super-eruptions, and nuclear war would all eject massive amounts of soot and aerosols into the atmosphere, with significant effects on global climate. The existence of such causal linkages is one reason why it is can be sensible to study multiple risks together.

Another commonality is that many methodological, conceptual, and cultural issues crop up across the range of global catastrophic risks. If our interest lies in such issues, it is often illuminating to study how they play out in different contexts. Conversely, some general insights – for example, into the biases of human risk cognition – can be applied to many different risks and used to improve our assessments across the board.

Beyond these theoretical commonalities, there are also pragmatic reasons for addressing global catastrophic risks as a single field. Attention is scarce. Mitigation is costly. To decide how to allocate effort and resources, we must make comparative judgements. If we treat risks singly, and never as part of an overall threat profile, we may become unduly fixated on the one or two dangers that happen to have captured the public or expert imagination of the day, while neglecting other risks that are more severe or more amenable to mitigation. Alternatively, we may fail to see that some precautionary policy, while effective in reducing the particular risk we are focusing on, would at the same time create new hazards and result in an increase in the overall level of risk. A broader view allows us to gain perspective and can thereby help us to set wiser priorities.

The immediate aim of this book is to offer an introduction to the range of global catastrophic risks facing humanity now or expected in the future, suitable for an educated interdisciplinary readership. There are several constituencies for the knowledge presented. Academics specializing in one of these risk areas will benefit from learning about the other risks. Professionals in insurance, finance, and business – although usually preoccupied with more limited and imminent challenges – will benefit from a wider view. Policy analysts, activists, and laypeople concerned with promoting responsible policies likewise stand to gain from learning about the state of the art in global risk studies. Finally, anyone who is worried or simply curious about what could go wrong in the modern world might find many of the following chapters intriguing. We hope that this volume will serve as a useful introduction to all of these audiences. Each of the chapters ends with some pointers to the literature for those who wish to delve deeper into a particular set of issues.

This volume also has a wider goal: to stimulate increased research, awareness, and informed public discussion about big risks and mitigation strategies. The existence of an interdisciplinary community of experts and laypeople knowledgeable about global catastrophic risks will, we believe, improve the odds that good solutions will be found and implemented to the great challenges of the twenty-first century.

1.2 Taxonomy and organization

Let us look more closely at what would, and would not, count as a global catastrophic risk. Recall that the damage must be serious, and the scale global. Given this, a catastrophe that caused 10,000 fatalities or 10 billion dollars worth of economic damage (e.g., a major earthquake) would not qualify as a global catastrophe. A catastrophe that caused 10 million fatalities or 10 trillion dollars worth of economic loss (e.g., an influenza pandemic) would count as a global catastrophe, even if some region of the world escaped unscathed. As for

disasters falling between these points, the definition is vague. The stipulation of a precise cut-off does not appear needful at this stage.

Global catastrophes have occurred many times in history, even if we only count disasters causing more than 10 million deaths. A very partial list of examples might include the An Shi Rebellion (756–763), the Taiping Rebellion (1851–1864), and the famine of the Great Leap Forward in China, the Black Death in Europe, the Spanish flu pandemic, the two world wars, the Nazi genocides, the famines in British India, Stalinist totalitarianism, the decimation of the native American population through smallpox and other diseases following the arrival of European colonizers, probably the Mongol conquests, perhaps Belgian Congo – innumerable others could be added to the list depending on how various misfortunes and chronic conditions are individuated and classified.

We can roughly characterize the severity of a risk by three variables: its *scope* (how many people – and other morally relevant beings – would be affected), its *intensity* (how badly these would be affected), and its *probability* (how likely the disaster is to occur, according to our best judgement, given currently available evidence). Using the first two of these variables, we can construct a qualitative diagram of different types of risk (Fig. 1.1). (The probability dimension could be displayed along a *z*-axis were this diagram three-dimensional.)

The scope of a risk can be *personal* (affecting only one person), *local*, *global* (affecting a large part of the human population), or *trans-generational* (affecting

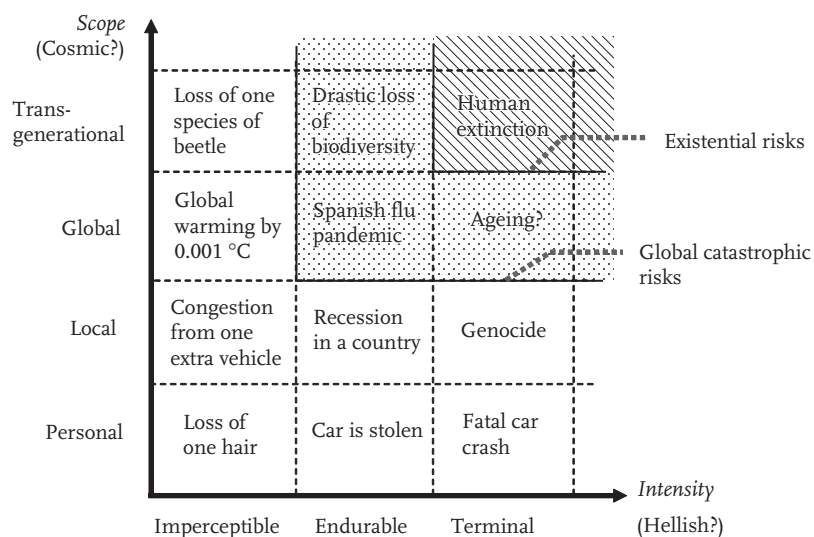


Fig. 1.1 Qualitative categories of risk. Global catastrophic risks are in the upper right part of the diagram. Existential risks form an especially severe subset of these.

not only the current world population but all generations that could come to exist in the future). The intensity of a risk can be classified as imperceptible (barely noticeable), endurable (causing significant harm but not destroying quality of life completely), or terminal (causing death or permanently and drastically reducing quality of life). In this taxonomy, global catastrophic risks occupy the four risks classes in the high-severity upper-right corner of the figure: a global catastrophic risk is of either global or trans-generational scope, and of either endurable or terminal intensity. In principle, as suggested in the figure, the axes can be extended to encompass conceptually possible risks that are even more extreme. In particular, trans-generational risks can contain a subclass of risks so destructive that their realization would not only affect or pre-empt future human generations, but would also destroy the potential of our future light cone of the universe to produce intelligent or self-aware beings (labelled ‘Cosmic’). On the other hand, according to many theories of value, there can be states of being that are even worse than non-existence or death (e.g., permanent and extreme forms of slavery or mind control), so it could, in principle, be possible to extend the x -axis to the right as well (see Fig. 1.1 labelled ‘Hellish’).

A subset of global catastrophic risks is *existential risks*. An existential risk is one that threatens to cause the extinction of Earth-originating intelligent life or to reduce its quality of life (compared to what would otherwise have been possible) permanently and drastically.¹ Existential risks share a number of features that mark them out as deserving of special consideration. For example, since it is not possible to recover from existential risks, we cannot allow even one existential disaster to happen; there would be no opportunity to learn from experience. Our approach to managing such risks must be proactive. How much worse an existential catastrophe would be than a non-existential global catastrophe depends very sensitively on controversial issues in value theory, in particular how much weight to give to the lives of possible future persons.² Furthermore, assessing existential risks raises distinctive methodological problems having to do with observation selection effects and the need to avoid anthropic bias. One of the motives for producing this book is to stimulate more serious study of existential risks. Rather than limiting our focus to existential risk, however, we thought it better to lay a broader foundation of systematic thinking about big risks in general.

¹ (Bostrom, 2002, p. 381).

² For many aggregative consequentialist ethical theories, including but not limited to total utilitarianism, it can be shown that the injunction to *maximize expected value!* can be simplified – for all practical purposes – to the injunction to *minimize existential risk!* (Bostrom, 2003, p. 439). (Note, however, that aggregative consequentialism is threatened by the problem of infinitarian paralysis [Bostrom, 2007, p. 730].)

We asked our contributors to assess global catastrophic risks not only as they presently exist but also as they might develop over time. The temporal dimension is essential for a full understanding of the nature of the challenges we face. To think about how to tackle the risks from nuclear terrorism and nuclear war, for instance, we must consider not only the probability that something will go wrong within the next year, but also about how the risks will change in the future and the factors – such as the extent of proliferation of relevant technology and fissile materials – that will influence this. Climate change from greenhouse gas emissions poses no significant globally catastrophic risk now or in the immediate future (on the timescale of several decades); the concern is about what effects these accumulating emissions might have over the course of many decades or even centuries. It can also be important to anticipate hypothetical risks which will arise if and when certain possible technological developments take place. The chapters on nanotechnology and artificial intelligence are examples of such prospective risk analysis.

In some cases, it can be important to study scenarios which are almost certainly physically impossible. The hypothetical risk from particle collider experiments is a case in point. It is very likely that these experiments have no potential, whatever, for causing global disasters. The objective risk is probably zero, as believed by most experts. But just how confident can we be that there is no objective risk? If we are not certain that there is no objective risk, then there is a risk at least in a subjective sense. Such subjective risks can be worthy of serious consideration, and we include them in our definition of global catastrophic risks.

The distinction between objective and subjective (epistemic) risk is often hard to make out. The possibility of an asteroid colliding with Earth looks like a clear-cut example of objective risk. But suppose that in fact no sizeable asteroid is on collision course with our planet within a certain, sufficiently large interval of time. We might then say that there is no objective risk of an asteroid-caused catastrophe within that interval of time. Of course, we will not know that this is so until we have mapped out the trajectories of all potentially threatening asteroids and are able to calculate all perturbations, often chaotic, of those trajectories. In the meantime, we must recognize a risk from asteroids even though the risk might be purely subjective, merely reflecting our present state of ignorance. An empty cave can be similarly subjectively unsafe if you are unsure about whether a lion resides in it; and it can be rational for you to avoid the cave if you reasonably judge that the expected harm of entry outweighs the expected benefit.

In the case of the asteroid threat, we have access to plenty of data that can help us quantify the risk. We can estimate the probability of a

catastrophic impact from statistics of past impacts (e.g., cratering data) and from observations sampling from the population of non-threatening asteroids. This particular risk, therefore, lends itself to rigorous scientific study, and the probability estimates we derive are fairly strongly constrained by hard evidence.³

For many other risks, we lack the data needed for rigorous statistical inference. We may also lack well-corroborated scientific models on which to base probability estimates. For example, there exists no rigorous scientific way of assigning a probability to the risk of a serious terrorist attack employing a biological warfare agent occurring within the next decade. Nor can we firmly establish that the risks of a global totalitarian regime arising before the end of the century are of a certain precise magnitude. It is inevitable that analyses of such risks will rely to a large extent on plausibility arguments, analogies, and subjective judgement.

Although more rigorous methods are to be preferred whenever they are available and applicable, it would be misplaced scientism to confine attention to those risks that are amenable to hard approaches.⁴ Such a strategy would lead to many risks being ignored, including many of the largest risks confronting humanity. It would also create a false dichotomy between two types of risks – the ‘scientific’ ones and the ‘speculative’ ones – where, in reality, there is a continuum of analytic tractability.

We have, therefore, opted to cast our net widely. Although our topic selection shows some skew towards smaller risks that have been subject to more scientific study, we do have a range of chapters that tackle potentially large but more speculative risks. The page count allocated to a risk should not, of course, be interpreted as a measure of how seriously we believe the risk ought to be regarded. In some cases, we have seen it fit to have a chapter devoted to a risk that turns out to be quite small, because learning that a particular risk is small can be useful, and the procedures used to arrive at the conclusion might serve as a template for future risk research. It goes without saying that the exact composition of a volume like this is also influenced by many contingencies

³ One can sometimes define something akin to objective physical probabilities (‘chances’) for deterministic systems, as is done, for example, in classical statistical mechanics, by assuming that the system is ergodic under a suitable course graining of its state space. But ergodicity is not necessary for there being strong scientific constraints on subjective probability assignments to uncertain events in deterministic systems. For example, if we have good statistics going back a long time showing that impacts occur on average once per thousand years, with no apparent trends or periodicity, then we have scientific reason – absent of more specific information – for assigning a probability of $\approx 0.1\%$ to an impact occurring within the next year, whether we think the underlying system dynamic is indeterministic, or chaotic, or something else.

⁴ Of course, when allocating research effort it is legitimate to take into account not just how important a problem is but also the likelihood that a solution can be found through research. The drunk who searches for his lost keys where the light is best is not necessarily irrational; and a scientist who succeeds in something relatively unimportant may achieve more good than one who fails in something important.

beyond the editors' control and that perforce it must leave out more than it includes.⁵

We have divided the book into four sections:

Part I: Background

Part II: Risks from Nature

Part III: Risks from Unintended Consequences

Part IV: Risks from Hostile Acts

This subdivision into three categories of risks is for convenience only, and the allocation of a risk to one of these categories is often fairly arbitrary. Take earthquakes which might seem to be paradigmatically a 'Risk from Nature'. Certainly, an earthquake is a natural event. It would happen even if we were not around. Earthquakes are governed by the forces of plate tectonics over which human beings currently have no control. Nevertheless, the risk posed by an earthquake is, to a very large extent, a matter of human construction. Where we erect our buildings and how we choose to construct them strongly influence what happens when an earthquake of a given magnitude occurs. If we all lived in tents, or in earthquake-proof buildings, or if we placed our cities far from fault lines and sea shores, earthquakes would do little damage. On closer inspection, we thus find that the earthquake risk is very much a joint venture between Nature and Man. Or take a paradigmatically anthropogenic hazard such as nuclear weapons. Again we soon discover that the risk is not as disconnected from uncontrollable forces of nature as might at first appear to be the case. If a nuclear bomb goes off, how much damage it causes will be significantly influenced by the weather. Wind, temperature, and precipitation will affect the fallout pattern and the likelihood that a fire storm will break out: factors that make a big difference to the number of fatalities generated by the blast. In addition, depending on how a risk is defined, it may also over time transition from one category to another. For instance, the risk of starvation might once have been primarily a Risk from Nature, when the main causal factors were draughts or fluctuations in local prey population; yet in the contemporary world, famines tend to be the consequences of market failures, wars, and social breakdowns, whence the risk is now at least as much one of Unintended Consequences or of Hostile Acts.

1.3 Part I: Background

The objective of this part of the book is to provide general context and methodological guidance for thinking systematically and critically about global catastrophic risks.

⁵ For example, the risk of large-scale conventional war is only covered in passing, yet would surely deserve its own chapter in a more ideally balanced page allocation.

We begin at the end, as it were, with Chapter 2 by Fred Adams discussing the long-term fate of our planet, our galaxy, and the Universe in general. In about 3.5 billion years, the growing luminosity of the sun will essentially have sterilized the Earth's biosphere, but the end of *complex* life on Earth is scheduled to come sooner, maybe 0.9–1.5 billion years from now. This is the default fate for life on our planet. One may hope that if humanity and complex technological civilization survives, it will long before then have learned to colonize space.

If some cataclysmic event were to destroy *Homo sapiens* and other higher organisms on Earth tomorrow, there does appear to be a window of opportunity of approximately one billion years for another intelligent species to evolve and take over where we left off. For comparison, it took approximately 1.2 billion years from the rise of sexual reproduction and simple multicellular organisms for the biosphere to evolve into its current state, and only a few million years for our species to evolve from its anthropoid ancestors. Of course, there is no guarantee that a rerun of evolution would produce anything like a human or a self-aware successor species.

If intelligent life does spread into space by harnessing the powers of technology, its lifespan could become extremely long. Yet eventually, the universe will wind down. The last stars will stop shining 100 trillion years from now. Later, matter itself will disintegrate into its basic constituents. By 10^{100} years from now even the largest black holes would have evaporated. Our present understanding of what will happen at this time scale and beyond is quite limited. The current best guess – but it is really no more than that – is that it is not just technologically difficult but physically impossible for intelligent information processing to continue beyond some finite time into the future. If so, extinction is not a question of whether, but when.

After this peek into the extremely remote future, it is instructive to turn around and take a brief peek at the distant past. Some past cataclysmic events have left traces in the geological record. There have been about 15 mass extinctions in the last 500 million years, and 5 of these eliminated more half of all species then inhabiting the Earth. Of particular note is the Permian – Triassic extinction event, which took place some 251.4 million years ago. This 'mother of all mass extinctions' eliminated more than 90% of all species and many entire phylogenetic families. It took upwards of 5 million years for biodiversity to recover.

Impacts from asteroids and comets, as well as massive volcano eruptions, have been implicated in many of the mass extinctions of the past. Other causes, such as variations in the intensity of solar illumination, may in some cases have exacerbated stresses. It appears that all mass extinctions have been mediated by atmospheric effects such as changes the atmosphere's composition or temperature. It is possible, however, that we owe our existence to mass extinctions. In particular, the comet that hit Earth 65 million years ago, which

is believed to have been responsible for the demise of the dinosaurs, might have been a *sine qua non* for the subsequent rise of *Homo sapiens* by clearing an ecological niche that could be occupied by large mammals, including our ancestors.

At least 99.9% of all species that have ever walked, crawled, flown, swum, or otherwise abided on Earth are extinct. Not all of these were eliminated in cataclysmic mass extinction events. Many succumbed in less spectacular doomsdays such as from competition by other species for the same ecological niche. Chapter 3 reviews the mechanisms of evolutionary change. Not so long ago, our own species co-existed with at least one other hominid species, the Neanderthals. It is believed that the lineages of *H. sapiens* and *H. neanderthalensis* diverged about 800,000 years ago. The Neanderthals manufactured and used composite tools such as handaxes. They did not reach extinction in Europe until 33,000 to 24,000 years ago, quite likely as a direct result of competition with *Homo sapiens*. Recently, the remains of what might have been another hominoid species, *Homo floresiensis* – nicknamed ‘the hobbit’ for its short stature – were discovered on an Indonesian island. *H. floresiensis* is believed to have survived until as recently as 12,000 years ago, although uncertainty remains about the interpretation of the finds. An important lesson of this chapter is that extinction of intelligent species *has* already happened on Earth, suggesting that it would be naïve to think it may not happen again.

From a naturalistic perspective, there is thus nothing abnormal about global cataclysms including species extinctions, although the characteristic time scales are typically large by human standards. James Hughes in Chapter 4 makes clear, however, the idea of cataclysmic endings often causes a peculiar set of cognitive tendencies to come into play, what he calls ‘the millennial, utopian, or apocalyptic psychocultural bundle, a characteristic dynamic of eschatological beliefs and behaviours’. The millennial impulse is pancultural. Hughes shows how it can be found in many guises and with many common tropes from Europe to India to China, across the last several thousand years. ‘We may aspire to a purely rational, technocratic analysis’, Hughes writes, ‘calmly balancing the likelihoods of futures without disease, hunger, work or death, on the one hand, against the likelihoods of worlds destroyed by war, plagues or asteroids, but few will be immune to millennial biases, positive or negative, fatalist or messianic’. Although these eschatological tropes can serve legitimate social needs and help to mobilize needed action, they easily become dysfunctional and contribute to social disengagement. Hughes argues that we need historically informed and vigilant self-interrogation to help us keep our focus on constructive efforts to address real challenges.

Even for an honest, truth-seeking, and well-intentioned investigator it is difficult to think and act rationally in regard to global catastrophic risks and existential risks. These are topics on which it seems especially

difficult to remain sensible. In Chapter 5, Eliezer Yudkowsky observes as follows:

Substantially larger numbers, such as 500 million deaths, and *especially* qualitatively different scenarios such as the extinction of the entire human species, seem to trigger a *different mode of thinking* – enter into a ‘separate magisterium’. People who would never dream of hurting a child hear of an existential risk, and say, ‘Well, maybe the human species doesn’t really deserve to survive’.

Fortunately, if we are ready to contend with our biases, we are not left entirely to our own devices. Over the last few decades, psychologists and economists have developed an extensive empirical literature on many of the common heuristics and biases that can be found in human cognition. Yudkowsky surveys this literature and applies its frequently disturbing findings to the domain of large-scale risks that is the subject matter of this book. His survey reviews the following effects: availability; hindsight bias; black swans; the conjunction fallacy; confirmation bias; anchoring, adjustment, and contamination; the affect heuristic; scope neglect; calibration and overconfidence; and bystander apathy. It behooves any sophisticated contributor in the area of global catastrophic risks and existential risks – whether scientist or policy advisor – to be familiar with each of these effects and we all ought to give some consideration to how they might be distorting our judgements.

Another kind of reasoning trap to be avoided is anthropic bias. Anthropic bias differs from the general cognitive biases reviewed by Yudkowsky; it is more theoretical in nature and it applies more narrowly to only certain specific kinds of inference. Anthropic bias arises when we overlook relevant observation selection effects. An observation selection effect occurs when our evidence has been ‘filtered’ by the precondition that a suitably positioned observer exists to have the evidence, in such a way that our observations are unrepresentatively sampled from the target domain. Failure to take observation effects into account correctly can result in serious errors in our probabilistic evaluation of some of the relevant hypotheses. Milan Ćirković, in Chapter 6, reviews some applications of observation selection theory that bear on global catastrophic risk and particularly existential risk. Some of these applications are fairly straightforward albeit not always obvious. For example, the tempting inference that certain classes of existential disaster must be highly improbable because they have never occurred in the history of our species or even in the history of life on Earth must be resisted. We are bound to find ourselves in one of those places and belonging to one of those intelligent species which have not yet been destroyed, whether planet or species-destroying disasters are common or rare: for the alternative possibility – that *our* planet has been destroyed or *our* species extinguished – is something that is unobservable for us, per definition. Other applications of anthropic reasoning – such as the Carter–Leslie Doomsday argument – are of disputed validity, especially

in their generalized forms, but nevertheless worth knowing about. In some applications, such as the simulation argument, surprising constraints are revealed on what we can coherently assume about humanity's future and our place in the world.

There are professional communities that deal with risk assessment on a daily basis. The subsequent two chapters present perspectives from the systems engineering discipline and the insurance industry, respectively.

In Chapter 7, Yacov Haimes outlines some flexible strategies for organizing our thinking about risk variables in complex systems engineering projects. What knowledge is needed to make good risk management decisions? Answering this question, Haimes says, 'mandates seeking the "truth" about the unknowable complex nature of emergent systems; it requires intellectually bias-free modellers and thinkers who are empowered to experiment with a multitude of modelling and simulation approaches and to collaborate for appropriate solutions'. Haimes argues that organizing the analysis around the measure of the expected value of risk can be too constraining. Decision makers often prefer a more fine-grained decomposition of risk that allows them to consider separately the probability of outcomes in different severity ranges, using what Haimes calls 'the partitioned multi-objective risk method'.

Chapter 8, by Peter Taylor, explores the connections between the insurance industry and global catastrophic risk. Insurance companies help individuals and organizations mitigate the financial consequences of risk, essentially by allowing risks to be traded and shared. Peter Taylor argues that the extent to which global catastrophic risks can be privately insured is severely limited for reasons having to do with both their scope and their type.

Although insurance and reinsurance companies have paid relatively scant attention to global catastrophic risks, they have accumulated plenty of experience with smaller risks. Some of the concepts and methods used can be applied to risks at any scale. Taylor highlights the importance of the concept of *uncertainty*. A particular stochastic model of phenomena in some domain (such as earthquakes) may entail a definite probability distribution over possible outcomes. However, in addition to the chanciness described by the model, we must recognize two further sources of uncertainty. There is usually uncertainty in the values of the parameters that we feed into the model. On top of that, there is uncertainty about whether the model we use does, in fact, correctly describe the phenomena in the target domain. These higher-level uncertainties are often impossible to analyse in a statistically rigorous way. Analysts who strive for objectivity and who are expected to avoid making 'un-scientific' assumptions that they cannot justify face a temptation to ignore these subjective uncertainties. But such scientism can lead to disastrous misjudgements. Taylor argues that the distortion is often greatest at the tail end of exceedance probability curves, leading to an underestimation of the risk of extreme events.

Taylor also reports on two recent survey studies of perceived risk. One of these, conducted by Swiss Re in 2005, asked executives of multinationals about which risks to their businesses' financials were of greatest concern to them. Computer-related risk was rated as the highest priority risk, followed by foreign trade, corporate governance, operational/facility, and liability risk. Natural disasters came in seventh place, and terrorism in tenth place. It appears that, as far as financial threats to individual corporations are concerned, global catastrophic risks take the backseat to more direct and narrowly focused business hazards. A similar exercise, but with broader scope, is carried out annually by the World Economic Forum. Its 2007 Global Risk report classified risks by likelihood and severity based on opinions solicited from business leaders, economists, and academics. Risks were evaluated with a 10-year time frame. Two risks were given a severity rating of 'more than 1 trillion USD', namely, asset price collapse (10–20%) and retrenchment from globalization (1–5%). When severity was measured in number of deaths rather than economic losses, the top three risks were pandemics, developing world disease, and interstate and civil war. (Unfortunately, several of the risks in this survey were poorly defined, making it hard to interpret the reported opinions – one moral here being that, if one wishes to assign probabilities to risks or rank them according to severity or likelihood, an essential first step is to present clear definitions of the risks that are to be evaluated.⁶)

The Background part of the book ends with a discussion by Richard Posner on some challenges for public policy in Chapter 9. Posner notes that governmental action to reduce global catastrophic risk is often impeded by the short decision horizons of politicians with their limited terms of office and the many competing demands on their attention. Furthermore, mitigation of global catastrophic risks is often costly and can create a free-rider problem. Smaller and poorer nations may drag their heels in the hope of taking a free ride on larger and richer countries. The more resourceful countries, in turn, may hold back because of reluctance to reward the free riders.

Posner also looks at several specific cases, including tsunamis, asteroid impacts, bioterrorism, accelerator experiments, and global warming, and considers some of the implications for public policy posed by these risks. Although rigorous cost–benefit analyses are not always possible, it is nevertheless important to attempt to quantify probabilities, potential harms, and the costs of different possible countermeasures, in order to determine priorities and optimal strategies for mitigation. Posner suggests that when

⁶ For example, the risk 'Chronic disease in the developed world' is defined as 'Obesity, diabetes and cardiovascular diseases become widespread; healthcare costs increase; resistant bacterial infections rise, sparking class-action suits and avoidance of hospitals'. By most standards, obesity, diabetes, and cardiovascular disease are *already* widespread. And by *how much* would healthcare costs have to increase to satisfy the criterion? It may be impossible to judge whether this definition was met even after the fact and with the benefit of hindsight.

a precise probability of some risk cannot be determined, it can sometimes be informative to consider – as a rough heuristic – the ‘implied probability’ suggested by current expenditures on mitigation efforts compared to the magnitude of harms that would result if a disaster materialized. For example, if we spend one million dollars per year to mitigate a risk which would create 1 billion dollars of damage, we may estimate that current policies implicitly assume that the annual risk of the disaster is of the order of 1/1000. If this implied probability seems too small, it might be a sign that we are not spending enough on mitigation.⁷ Posner maintains that the world is, indeed, under-investing in mitigation of several global catastrophic risks.

1.4 Part II: Risks from nature

Volcanic eruptions in recent historical times have had measurable effects on global climate, causing global cooling by a few tenths of one degree, the effect lasting perhaps a year. But as Michael Rampino explains in Chapter 10, these eruptions pale in comparison to the largest recorded eruptions. Approximately 75,000 years ago, a volcano erupted in Toba, Indonesia, spewing vast volumes of fine ash and aerosols into the atmosphere, with effects comparable to nuclear-winter scenarios. Land temperatures globally dropped by 5–15°C, and ocean-surface cooling of $\approx 2\text{--}6^\circ\text{C}$ might have extended over several years. The persistence of significant soot in the atmosphere for 1–3 years might have led to a cooling of the climate lasting for decades (because of climate feedbacks such as increased snow cover and sea ice causing more of the sun’s radiation to be reflected back into space). The human population appears to have gone through a bottleneck at this time, according to some estimates dropping as low as ≈ 500 reproducing females in a world population of approximately 4000 individuals. On the Toba catastrophe theory, the population decline was caused by the super-eruption, and the human species was teetering on the brink of extinction. This is perhaps the worst disaster that has ever befallen the human species, at least if severity is measured by how close to terminal was the outcome.

More than 20 super-eruption sites for the last 2 million years have been identified. This would suggest that, on average, a super-eruption occurs at least once every 50,000 years. However, there may well have been additional super-eruptions that have not yet been identified in the geological record.

⁷ This heuristic is only meant to be a first stab at the problem. It is obviously not generally valid. For example, if one million dollars is sufficient to take all the possible precautions, there is no reason to spend more on the risk even if we think that its probability is much greater than 1/1000. A more careful analysis would consider the marginal returns on investment in risk reduction.

The global damage from super-volcanism would come chiefly from its climatic effects. The volcanic winter that would follow such an eruption would cause a drop in agricultural productivity which could lead to mass starvation and consequent social upheavals. Rampino's analysis of the impacts of super-volcanism is also relevant to the risks of nuclear war and asteroid or meteor impacts. Each of these would involve soot and aerosols being injected into the atmosphere, cooling the Earth's climate.

Although we have no way of preventing a super-eruption, there are precautions that we could take to mitigate its impacts. At present, a global stockpile equivalent to a 2-month supply of grain exists. In a super-volcanic catastrophe, growing seasons might be curtailed for several years. A larger stockpile of grain and other foodstuffs, while expensive to maintain, would provide a buffer for a range of catastrophe scenarios involving temporary reductions in world agricultural productivity.

The hazard from comets and meteors is perhaps the best understood of all global catastrophic risks (which is not to deny that significant uncertainties remain). Chapter 11, by William Napier, explains some of the science behind the impact hazards: where comets and asteroids come from, how frequently impacts occur, and what the effects of an impact would be. To produce a civilization-disrupting event, an impactor would need a diameter of at least 1 or 2 km. A 10-km impactor would, it appears, have a good chance of causing the extinction of the human species. But even sub-kilometre impactors could produce damage reaching the level of global catastrophe, depending on their composition, velocity, angle, and impact site.

Napier estimates that 'the per capita impact hazard is at the level associated with the hazards of air travel and the like'. However, funding for mitigation is meager compared to funding for air safety. The main effort currently underway to address the impact hazard is the *Spaceguard* project, which receives about 4 million dollars per annum from NASA besides in-kind and voluntary contributions from others. *Spaceguard* aims to find 90% of near-Earth asteroids larger than 1 km by the end of 2008. Asteroids constitute the largest portion of the threat from near-Earth objects (and are easier to detect than comets) so when the project is completed, the subjective probability of a large impact will have been reduced considerably – unless, of course, it were discovered that some asteroid has a date with our planet in the near future, in which case the probability would soar.

Some preliminary study has been done of how a potential impactor could be deflected. Given sufficient advance warning, it appears that the space technology needed to divert an asteroid could be developed. The cost of producing an effective asteroid defence would be much greater than the cost of searching for potential impactors. However, if a civilization-destroying wrecking ball were found to be swinging towards the Earth, virtually any expense would be justified to avert it before it struck.

Asteroids and comets are not the only potential global catastrophic threats from space. Other cosmic hazards include global climatic change from fluctuations in solar activity, and very large fluxes from radiation and cosmic rays from supernova explosions or gamma ray bursts. These risks are examined in Chapter 12 by Arnon Dar. The findings on these risks are favourable: the risks appear to be very small. No particular response seems indicated at the present time beyond continuation of basic research.⁸

1.5 Part III: Risks from unintended consequences

We have already encountered climate change – in the form of sudden global cooling – as a destructive modality of super-eruptions and large impacts (as well as possible consequence of large-scale nuclear war, to be discussed later). Yet it is the risk of gradual global warming brought about by greenhouse gas emissions that has most strongly captured the public imagination in recent years. Anthropogenic climate change has become the poster child of global threats. Global warming commandeers a disproportionate fraction of the attention given to global risks.

Carbon dioxide and other greenhouse gases are accumulating in the atmosphere, where they are expected to cause a warming of Earth's climate and a concomitant rise in seawater levels. The most recent report by the United Nations' Intergovernmental Panel on Climate Change (IPCC), which represents the most authoritative assessment of current scientific opinion, attempts to estimate the increase in global mean temperature that would be expected by the end of this century under the assumption that no efforts at mitigation are made. The final estimate is fraught with uncertainty because of uncertainty about what the default rate of emissions of greenhouse gases will be over the century, uncertainty about the climate sensitivity parameter, and uncertainty about other factors. The IPCC, therefore, expresses its assessment in terms of six different climate scenarios based on different models and different assumptions. The 'low' model predicts a mean global warming of +1.8°C (uncertainty range 1.1–2.9°C); the 'high' model predicts warming by +4.0°C (2.4–6.4°C). Estimated sea level rise predicted by the two most extreme scenarios of the six considered is 18–38 cm, and 26–59 cm, respectively.

Chapter 13, by David Frame and Myles Allen, summarizes some of the basic science behind climate modelling, with particular attention to the low-probability high-impact scenarios that are most relevant to the focus of this book. It is, arguably, this range of extreme scenarios that gives the greatest

⁸ A comprehensive review of space hazards would also consider scenarios involving contact with intelligent extraterrestrial species or contamination from hypothetical extraterrestrial microorganisms; however, these risks are outside the scope of Chapter 12.

cause for concern. Although their likelihood seems very low, considerable uncertainty still pervades our understanding of various possible feedbacks that might be triggered by the expected climate forcing (recalling Peter Taylor's point, referred to earlier, about the importance of taking parameter and model uncertainty into account). David Frame and Myles Allen also discuss mitigation policy, highlighting the difficulties of setting appropriate mitigation goals given the uncertainties about what levels of cumulative emissions would constitute 'dangerous anthropogenic interference' in the climate system.

Edwin Kilbourne reviews some historically important pandemics in Chapter 14, including the distinctive characteristics of their associated pathogens, and discusses the factors that will determine the extent and consequences of future outbreaks.

Infectious disease has exacted an enormous toll of suffering and death on the human species throughout history and continues to do so today. Deaths from infectious disease currently account for approximately 25% of all deaths worldwide. This amounts to approximately 15 million deaths per year. About 75% of these deaths occur in Southeast Asia and sub-Saharan Africa. The top five causes of death due to infectious disease are upper respiratory infection (3.9 million deaths), HIV/AIDS (2.9 million), diarrhoeal disease (1.8 million), tuberculosis (1.7 million), and malaria (1.3 million).

Pandemic disease is indisputably one of the biggest global catastrophic risks facing the world today, but it is not always accorded its due recognition. For example, in most people's mental representation of the world, the influenza pandemic of 1918–1919 is almost completely overshadowed by the concomitant World War I. Yet although the WWI is estimated to have directly caused about 10 million military and 9 million civilian fatalities, the Spanish flu is believed to have killed at least 20–50 million people. The relatively low 'dread factor' associated with this pandemic might be partly due to the fact that only approximately 2–3% of those who got sick died from the disease. (The total death count is vast because a large percentage of the world population was infected.)

In addition to fighting the major infectious diseases currently plaguing the world, it is vital to remain alert to emerging new diseases with pandemic potential, such as SARS, bird flu, and drug-resistant tuberculosis. As the World Health Organization and its network of collaborating laboratories and local governments have demonstrated repeatedly, decisive early action can sometimes nip an emerging pandemic in the bud, possibly saving the lives of millions.

We have chosen to label pandemics a 'risk from unintended consequences' even though most infectious diseases (exempting the potential of genetically engineered bioweapons) in some sense arise from nature. Our rationale is that the evolution as well as the spread of pathogens is highly dependent on human civilization. The worldwide spread of germs became possible only after all the

inhabited continents were connected by travel routes. By now, globalization in the form of travel and trade has reached such an extent that a highly contagious disease could spread to virtually all parts of the world within a matter of days or weeks. Kilbourne also draws attention to another aspect of globalization as a factor increasing pandemic risk: homogenization of peoples, practices, and cultures. The more the human population comes to resemble a single homogeneous niche, the greater the potential for a single pathogen to saturate it quickly. Kilbourne mentions the ‘one rotten apple syndrome’, resulting from the mass production of food and behavioural fads:

If one contaminated item, apple, egg or most recently spinach leaf carries a billion bacteria – not an unreasonable estimate – and it enters a pool of cake mix constituents then packaged and sent to millions of customers nationwide, a bewildering epidemic may ensue.

Conversely, cultural as well as genetic diversity reduces the likelihood that any single pattern will be adopted universally before it is discovered to be dangerous – whether the pattern be virus RNA, a dangerous new chemical or material, or a stifling ideology.

By contrast to pandemics, artificial intelligence (AI) is not an ongoing or imminent global catastrophic risk. Nor is it as uncontroversially a serious cause for concern. However, from a long-term perspective, the development of general artificial intelligence exceeding that of the human brain can be seen as one of the main challenges to the future of humanity (arguably, even as *the* main challenge). At the same time, the successful deployment of friendly superintelligence could obviate many of the other risks facing humanity. The title of Chapter 15, ‘Artificial intelligence as a positive and negative factor in global risk’, reflects this ambivalent potential.

As Eliezer Yudkowsky notes, the prospect of superintelligent machines is a difficult topic to analyse and discuss. Appropriately, therefore, he devotes a substantial part of his chapter to clearing common misconceptions and barriers to understanding. Having done so, he proceeds to give an argument for giving serious consideration to the possibility that radical superintelligence could erupt very suddenly – a scenario that is sometimes referred to as the ‘Singularity hypothesis’. Claims about the steepness of the transition must be distinguished from claims about the timing of its onset. One could believe, for example, that it will be a long time before computers are able to match the general reasoning abilities of an average human being, but that once that happens, it will only take a short time for computers to attain radically superhuman levels.

Yudkowsky proposes that we conceive of a superintelligence as an enormously powerful optimization process: ‘a system which hits small targets in large search spaces to produce coherent real-world effects’. The superintelligence will be able to manipulate the world (including human beings) in such a way as to achieve its goals, whatever those goals might be.

To avert disaster, it would be necessary to ensure that the superintelligence is endowed with a 'Friendly' goal system: that is, one that aligns the system's goals with genuine human values.

Given this set-up, Yudkowsky identifies two different ways in which we could fail to build Friendliness into our AI: philosophical failure and technical failure. The warning against philosophical failure is basically that we should be careful what we wish for because we might get it. We might designate a target for the AI which at first sight seems like a nice outcome but which in fact is radically misguided or morally worthless. The warning against technical failure is that we might fail to get what we wish for, because of faulty implementation of the goal system or unintended consequences of the way the target representation was specified. Yudkowsky regards both of these possible failure modes as very serious existential risks and concludes that it is imperative that we figure out how to build Friendliness into a superintelligence before we figure out how to build a superintelligence.

Chapter 16 discusses the possibility that the experiments that physicists carry out in particle accelerators might pose an existential risk. Concerns about such risks prompted the director of the Brookhaven Relativistic Heavy Ion Collider to commission an official report in 2000. Concerns have since resurfaced with the construction of more powerful accelerators such as CERN's Large Hadron Collider. Following the Brookhaven report, Frank Wilczek distinguishes three catastrophe scenarios:

1. Formation of tiny black holes that could start accreting surrounding matter, eventually swallowing up the entire planet.
2. Formation of negatively charged stable strangelets which could catalyse the conversion of all the ordinary matter on our planet into strange matter.
3. Initiation of a phase transition of the vacuum state, which would propagate outward in all directions at near light speed and destroy not only our planet but the entire accessible part of the universe.

Wilczek argues that these scenarios are exceedingly unlikely on various theoretical grounds. In addition, there is a more general argument that these scenarios are extremely improbable which depends less on arcane theory. Cosmic rays often have energies far greater than those that will be attained in any of the planned accelerators. Such rays have been bombarding the Earth's atmosphere (and the moon and other astronomical objects) for billions of years without a single catastrophic effect having been observed. Assuming that collisions in particle accelerators do not differ in any unknown relevant respect from those that occur in the wild, we can be very confident in the safety of our accelerators.

By everyone's reckoning, it is highly improbable that particle accelerator experiments will cause an existential disaster. The question is *how* improbable? And what would constitute an 'acceptable' probability of an existential disaster?

In assessing the probability, we must consider not only how unlikely the outcome seems given our best current models but also the possibility that our best models and calculations might be flawed in some as-yet unrealized way. In doing so we must guard against overconfidence bias (compare Chapter 5 on biases). Unless we ourselves are technically expert, we must also take into account the possibility that the experts on whose judgements we rely might be consciously or unconsciously biased.⁹ For example, the physicists who possess the expertise needed to assess the risks from particle physics experiments are part of a professional community that has a direct stake in the experiments going forward. A layperson might worry that the incentives faced by the experts could lead them to err on the side of downplaying the risks.¹⁰ Alternatively, some experts might be tempted by the media attention they could get by playing up the risks. The issue of how much and in which circumstances to trust risk estimates by experts is an important one, and it arises quite generally with regard to many of the risks covered in this book.

Chapter 17 (by Robin Hanson) from Part III on Risks from unintended consequences focuses on social collapse as a devastation multiplier of other catastrophes. Hanson writes as follows:

The main reason to be careful when you walk up a flight of stairs is not that you might slip and have to retrace one step, but rather that the first slip might cause a second slip, and so on until you fall dozens of steps and break your neck. Similarly we are concerned about the sorts of catastrophes explored in this book not only because of their terrible direct effects, but also because they may induce an even more damaging collapse of our economic and social systems.

This argument does not apply to some of the risks discussed so far, such as those from particle accelerators or the risks from superintelligence as envisaged by Yudkowsky. In those cases, we may be either completely safe or altogether doomed, with little probability of intermediary outcomes. But for many other types of risk – such as windstorms, tornados, earthquakes, floods, forest fires, terrorist attacks, plagues, and wars – a wide range of outcomes are possible, and the potential for social disruption or even social collapse constitutes a major part of the overall hazard. Hanson notes that many of these risks appear to follow a power law distribution. Depending on the characteristic exponent of such a power law distribution, most of the damage expected from a given

⁹ Even if we ourselves are expert, we must still be alert to unconscious biases that may influence our judgment (e.g., anthropic biases, see Chapter 6).

¹⁰ If experts anticipate that the public will not quite trust their reassurances, they might be led to try to sound even more reassuring than they would have if they had believed that the public would accept their claims at face value. The public, in turn, might respond by discounting the experts' verdicts even more, leading the experts to be even more wary of fuelling alarmist overreactions. In the end, experts might be reluctant to acknowledge any risk at all for fear of a triggering a hysterical public overreaction. Effective risk communication is a tricky business, and the trust that it requires can be hard to gain and easy to lose.

type of risk may consist either of frequent small disturbances or of rare large catastrophes. Car accidents, for example, have a large exponent, reflecting the fact that most traffic deaths occur in numerous small accidents involving one or two vehicles. Wars and plagues, by contrast, appear to have small exponents, meaning that most of the expected damage occurs in very rare but very large conflicts and pandemics.

After giving a thumbnail sketch of economic growth theory, Hanson considers an extreme opposite of economic growth: sudden reduction in productivity brought about by escalating destruction of social capital and coordination. For example, 'a judge who would not normally consider taking a bribe may do so when his life is at stake, allowing others to expect to get away with theft more easily, which leads still others to avoid making investments that might be stolen, and so on. Also, people may be reluctant to trust bank accounts or even paper money, preventing those institutions from functioning.' The productivity of the world economy depends both on scale and on many different forms of capital which must be delicately coordinated. We should be concerned that a relatively small disturbance (or combination of disturbances) to some vulnerable part of this system could cause a far-reaching unraveling of the institutions and expectations upon which the global economy depends.

Hanson also offers a suggestion for how we might convert some existential risks into non-existential risks. He proposes that we consider the construction of one or more continuously inhabited refuges – located, perhaps, in a deep mineshaft, and well-stocked with supplies – which could preserve a small but sufficient group of people to repopulate a post-apocalyptic world. It would obviously be preferable to prevent altogether catastrophes of a severity that would make humanity's survival dependent on such modern-day 'Noah's arks'; nevertheless, it might be worth exploring whether some variation of this proposal might be a cost-effective way of somewhat decreasing the probability of human extinction from a range of potential causes.¹¹

1.6 Part IV: Risks from hostile acts

The spectre of nuclear Armageddon, which so haunted the public imagination during the Cold War era, has apparently entered semi-retirement. The number of nuclear weapons in the world has been reduced to half, from a Cold War high of 65,000 in 1986 to approximately 26,000 in 2007, with approximately

¹¹ Somewhat analogously, we could prevent much permanent loss of biodiversity by moving more aggressively to preserve genetic material from endangered species in biobanks. The Norwegian government has recently opened a seed bank on a remote island in the arctic archipelago of Svalbard. The vault, which is dug into a mountain and protected by steel-reinforced concrete walls one metre thick, will preserve germplasm of important agricultural and wild plants.

96% of these weapons held by the United States and Russia. Relationships between these two nations are not as bad as they once were. New scares such as environmental problems and terrorism compete effectively for media attention. Changing winds in horror-fashion aside, however, and as Chapter 18 makes it clear, nuclear war remains a very serious threat.

There are several possibilities. One is that relations between the United States and Russia might again worsen to the point where a crisis could trigger a nuclear war. Future arms races could lead to arsenals even larger than those of the past. The world's supply of plutonium has been increasing steadily to about 2000 tons – about 10 times as much as remains tied up in warheads – and more could be produced. Some studies suggest that in an all-out war involving most of the weapons in the current US and Russian arsenals, 35–77% of the US population (105–230 million people) and 20–40% of the Russian population (28–56 million people) would be killed. Delayed and indirect effects – such as economic collapse and a possible nuclear winter – could make the final death toll far greater.

Another possibility is that nuclear war might erupt between nuclear powers other than the old Cold War rivals, a risk that is growing as more nations join the nuclear club, especially nations that are embroiled in volatile regional conflicts, such as India and Pakistan, North Korea, and Israel, perhaps to be joined by Iran or others. One concern is that the more nations get the bomb, the harder it might be to prevent further proliferation. The technology and know-how would become more widely disseminated, lowering the technical barriers, and nations that initially chose to forego nuclear weapons might feel compelled to rethink their decision and to follow suit if they see their neighbours start down the nuclear path.

A third possibility is that global nuclear war could be started by mistake. According to Joseph Cirincione, this almost happened in January 1995:

Russian military officials mistook a Norwegian weather rocket for a US submarine-launched ballistic missile. Boris Yelstin became the first Russian president to ever have the 'nuclear suitcase' open in front of him. He had just a few minutes to decide if he should push the button that would launch a barrage of nuclear missiles. Thankfully, he concluded that his radars were in error. The suitcase was closed.

Several other incidents have been reported in which the world, allegedly, was teetering on the brink of nuclear holocaust. At one point during the Cuban missile crises, for example, President Kennedy reportedly estimated the probability of a nuclear war between the United States and the USSR to be 'somewhere between one out of three and even'.

To reduce the risks, Cirincione argues, we must work to resolve regional conflicts, support and strengthen the Nuclear Non-proliferation Treaty – one of the most successful security pacts in history – and move towards the abolition of nuclear weapons.

William Potter and Gary Ackerman offer a detailed look at the risks of nuclear terrorism in Chapter 19. Such terrorism could take various forms:

- Dispersal of radioactive material by conventional explosives ('dirty bomb')
- Sabotage of nuclear facilities
- Acquisition of fissile material leading to the fabrication and detonation of a crude nuclear bomb ('improvised nuclear device')
- Acquisition and detonation of an intact nuclear weapon
- The use of some means to trick a nuclear state into launching a nuclear strike.

Potter and Ackerman focus on 'high consequence' nuclear terrorism, which they construe as those involving the last three alternatives from the above list. The authors analyse the demand and supply side of nuclear terrorism, the consequences of a nuclear terrorist attack, the future shape of the threat, and conclude with policy recommendations.

To date, no non-state actor is believed to have gained possession of a fission weapon:

There is no credible evidence that either al Qaeda or Aum Shinrikyo were able to exploit their high motivations, substantial financial resources, demonstrated organizational skills, far-flung network of followers, and relative security in a friendly or tolerant host country to move very far down the path toward acquiring a nuclear weapons capability. As best one can tell from the limited information available in public sources, among the obstacles that proved most difficult for them to overcome was access to the fissile material needed . . .

Despite this track record, however, many experts remain concerned. Graham Allison, author of one of the most widely cited works on the subject, offers a standing bet of 51 to 49 odds that 'barring radical new anti-proliferation steps' there will be a terrorist nuclear strike within the next 10 years. Other experts seem to place the odds much lower, but have apparently not taken up Allison's offer.

There is wide recognition of the importance of prevention nuclear terrorism, and in particular of the need to prevent fissile material from falling into the wrong hands. In 2002, the G-8 Global Partnership set a target of 20 billion dollars to be committed over a 10-year period for the purpose of preventing terrorists from acquiring weapons and materials of mass destruction. What Potter and Ackerman consider most lacking, however, is the sustained high-level leadership needed to transform rhetoric into effective implementation.

In Chapter 20, Christopher Chyba and Ali Nouri review issues related to biotechnology and biosecurity. While in some ways paralleling nuclear risks – biological as well as nuclear technology can be used to build weapons of mass destruction – there are also important divergences. One difference is that biological weapons can be developed in small, easily concealed facilities and

require no unusual raw materials for their manufacture. Another difference is that an infectious biological agent can spread far beyond the site of its original release, potentially across the entire world.

Biosecurity threats fall into several categories, including naturally occurring diseases, illicit state biological weapons programmes, non-state actors and bio-hackers, and laboratory accidents or other inadvertent release of disease agents. It is worth bearing in mind that the number of people who have died in recent years from threats in the first of these categories (naturally occurring diseases) is six or seven orders of magnitudes larger than the number of fatalities from the other three categories combined. Yet biotechnology does contain brewing threats which look set to expand dramatically over the coming years as capabilities advance and proliferate. Consider the following sample of recent developments:

- A group of Australian researchers, looking for ways of controlling the country's rabbit population, added the gene for interleukin-4 to a mousepox virus, hoping thereby to render the animals sterile. Unexpectedly, the virus inhibited the host's immune system and all the animals died, including individuals who had previously been vaccinated. Follow-up work by another group produced a version of the virus that was 100% lethal in vaccinated mice despite the antiviral medication given to the animals.
- The polio virus has been synthesized from readily purchased chemical supplies. When this was first done, it required a protracted cutting-edge research project. Since then, the time needed to synthesize a virus genome comparable in size to the polio virus has been reduced to weeks. The virus that caused the Spanish flu pandemic, which was previously extinct, has also been resynthesized and now exists in laboratories in the United States and in Canada.
- The technology to alter the properties of viruses and other microorganisms is advancing at a rapid pace. The recently developed method of RNA interference provides researchers with a ready means of turning off selected genes in humans and other organisms. 'Synthetic biology' is being established as new field, whose goal is to enable the creation of small biological devices and ultimately new types of microbes.

Reading this list, while bearing in mind that the complete genomes from hundreds of bacteria, fungi, viruses – including Ebola, Marburg, smallpox, and the 1918 Spanish influenza virus – have been sequenced and deposited in a public online database, it is not difficult to concoct in one's imagination frightening possibilities. The technological barriers to the production of super bugs are being steadily lowered even as the biotechnological know-how and equipment diffuse ever more widely.

The dual-use nature of the necessary equipment and expertise, and the fact that facilities could be small and easily concealed, pose difficult challenges for would-be regulators. For any regulatory regime to work, it would also have to strike a difficult balance between prevention of abuses and enablement of research needed to develop treatments and diagnostics (or to obtain other medical or economic benefits). Chyba and Nouri discuss several strategies for promoting biosecurity, including automated review of gene sequences submitted for DNA-synthesizing at centralized facilities. It is likely that biosecurity will grow in importance and that a multipronged approach will be needed to address the dangers from designer pathogens.

Chris Phoenix and Mike Treder (Chapter 21) discuss nanotechnology as a source of global catastrophic risks. They distinguish between ‘nanoscale technologies’, of which many exist today and many more are in development, and ‘molecular manufacturing’, which remains a hypothetical future technology (often associated with the person who first envisaged it in detail, K. Eric Drexler). Nanoscale technologies, they argue, appear to pose no new global catastrophic risks, although such technologies could in some cases either augment or help mitigate some of the other risks considered in this volume. Phoenix and Treder consequently devote the bulk of their chapter to considering the capabilities and threats from molecular manufacturing. As with superintelligence, the *present* risk is virtually zero since the technology in question does not yet exist; yet the future risk could be extremely severe.

Molecular nanotechnology would greatly expand control over the structure of matter. Molecular machine systems would enable fast and inexpensive manufacture of microscopic and macroscopic objects built to atomic precision. Such production systems would contain millions of microscopic assembly tools. Working in parallel, these would build objects by adding molecules to a workpiece through positionally controlled chemical reactions. The range of structures that could be built with such technology greatly exceeds that accessible to the biological molecular assemblers (such as ribosome) that exist in nature. Among the things that a nanofactory could build: another nanofactory. A sample of potential applications:

- microscopic nanobots for medical use
- vastly faster computers
- very light and strong diamondoid materials
- new processes for removing pollutants from the environment
- desktop manufacturing plants which can automatically produce a wide range of atomically precise structures from downloadable blueprints
- inexpensive solar collectors
- greatly improved space technology

- mass-produced sensors of many kinds
- weapons, both inexpensively mass-produced and improved conventional weapons, and new kinds of weapons that cannot be built without molecular nanotechnology.

A technology this powerful and versatile could be used for an indefinite number of purposes, both benign and malign.

Phoenix and Treder review a number of global catastrophic risks that could arise with such an advanced manufacturing technology, including war, social and economic disruption, destructive forms of global governance, radical intelligence enhancement, environmental degradation, and ‘ecophagy’ (small nanobots replicating uncontrollably in the natural environment, consuming or destroying the Earth’s biosphere). In conclusion, they offer the following rather alarming assessment:

In the absence of some type of preventive or protective force, the power of molecular manufacturing products could allow a large number of actors of varying types – including individuals, groups, corporations, and nations – to obtain sufficient capability to destroy all unprotected humans. The likelihood of at least one powerful actor being insane is not small. The likelihood that devastating weapons will be built and released accidentally (possibly through overly sensitive automated systems) is also considerable. Finally, the likelihood of a conflict between two [powers capable of unleashing a mutually assured destruction scenario] escalating until one feels compelled to exercise a doomsday option is also non-zero. This indicates that unless adequate defences can be prepared against weapons intended to be ultimately destructive – a point that urgently needs research – the number of actors trying to possess such weapons must be minimized.

The last chapter of the book, authored by Bryan Caplan, addresses totalitarianism as a global catastrophic risk. The totalitarian governments of Nazi Germany, Soviet Russia, and Maoist China were responsible for tens of millions of deaths in the last century. Compared to a risk like that of asteroid impacts, totalitarianism as a global risk is harder to study in an unbiased manner, and a cross-ideological consensus about how this risk is best to be mitigated is likely to be more elusive. Yet the risks from oppressive forms of government, including totalitarian regimes, must not be ignored. Oppression has been one of the major recurring banes of human development throughout history, it largely remains so today, and it is one to which the humanity remains vulnerable.

As Caplan notes, in addition to being a misfortune in itself, totalitarianism can also amplify other risks. People in totalitarian regimes are often afraid to publish bad news, and the leadership of such regimes is often insulated from criticism and dissenting views. This can make such regimes more likely to overlook looming dangers and to commit serious policy errors (even as evaluated from the standpoint of the self-interest of the rulers). However, as

Caplan notes further, for some types of risk, totalitarian regimes might actually possess an advantage compared to more open and diverse societies. For goals that can be achieved by brute force and massive mobilization of resources, totalitarian methods have often proven effective.

Caplan analyses two factors which he claims have historically limited the durability of totalitarian regimes. The first of these is the problem of succession. A strong leader might maintain a tight grip on power for as long as he lives, but the party faction he represents often stumbles when it comes to appointing a successor that will preserve the status quo, allowing a closet reformer – a sheep in wolf's clothing – to gain the leadership position after a tyrant's death. The other factor is the existence of non-totalitarian countries elsewhere in the world. These provide a vivid illustration to the people living under totalitarianism that things could be much better than they are, fuelling dissatisfaction and unrest. To counter this, leaders might curtail contacts with the external world, creating a 'hermit kingdom' such as Communist Albania or present-day North Korea. However, some information is bound to leak in. Furthermore, if the isolation is too complete, over a period of time, the country is likely to fall far behind economically and militarily, making itself vulnerable to invasion or externally imposed regime change.

It is possible that the vulnerability presented by these two Achilles heels of totalitarianism could be reduced by future developments. Technological advances could help solve the problem of succession. Brain scans might one day be used to screen out closet sceptics within the party. Other forms of novel surveillance technologies could also make it easier to control population. New psychiatric drugs might be developed that could increase docility without noticeably reducing productivity. Life-extension medicine might prolong the lifespan of the leader so that the problem of succession comes up less frequently. As for the existence of non-totalitarian outsiders, Caplan worries about the possible emergence of a world government. Such a government, even if it started out democratic, might at some point degenerate into totalitarianism; and a worldwide totalitarian regime could then have great staying power given its lack of external competitors and alien exemplars of the benefits of political freedom.

To have a productive discussion about matters such as these, it is important to recognize the distinction between two very different stances: 'here a valid consideration in favour of some position *X*' versus '*X* is all-things-considered the position to be adopted'. For instance, as Caplan notes:

If people lived forever, stable totalitarianism would be a little more likely to emerge, but it would be madness to force everyone to die of old age in order to avert a small risk of being murdered by the secret police in a thousand years.

Likewise, it is possible to favour the strengthening of certain new forms global governance while also recognizing as a legitimate concern the danger of global totalitarianism to which Caplan draws our attention.

1.7 Conclusions and future directions

The most likely global catastrophic risks all seem to arise from human activities, especially industrial civilization and advanced technologies. This is not necessarily an indictment of industry or technology, for these factors deserve much of the credit for creating the values that are now at risk – including most of the people living on the planet today, there being perhaps 30 times more of us than could have been sustained with primitive agricultural methods, and hundreds of times more than could have lived as hunter-gatherers. Moreover, although new global catastrophic risks have been created, many smaller-scale risks have been drastically reduced in many parts of the world, thanks to modern technological society. Local and personal disasters – such as starvation, thirst, predation, disease, and small-scale violence – have historically claimed many more lives than have global cataclysms. The reduction of the aggregate of these smaller-scale hazards may outweigh an increase in global catastrophic risks. To the (incomplete) extent that true risk levels are reflected in actuarial statistics, the world is a safer place than it has ever been: world life expectancy is now 64 years, up from 50 in the early twentieth century, 33 in Medieval Britain, and an estimated 18 years during the Bronze Age. Global catastrophic risks are, by definition, the largest in terms of *scope* but not necessarily in terms of their expected severity (probability \times harm). Furthermore, technology and complex social organizations offer many important tools for managing the remaining risks. Nevertheless, it is important to recognize that the biggest global catastrophic risks we face today are not purely external; they are, instead, tightly wound up with the direct and indirect, the foreseen and unforeseen, consequences of our own actions.

One major current global catastrophic risk is infectious pandemic disease. As noted earlier, infectious disease causes approximately 15 million deaths per year, of which 75% occur in Southeast Asia and sub-Saharan Africa. These dismal statistics pose a challenge to the classification of pandemic disease as a global catastrophic risk. One could argue that infectious disease is not so much a *risk* as an *ongoing global catastrophe*. Even on a more fine-grained individuation of the hazard, based on specific infectious agents, at least some of the currently occurring pandemics (such as HIV/AIDS, which causes nearly 3 million deaths annually) would presumably qualify as global catastrophes. By similar reckoning, one could argue that cardiovascular disease (responsible for approximately 30% of world mortality, or 18 million deaths per year) and cancer (8 million deaths) are also ongoing global catastrophes. It would be perverse if the study of possible catastrophes that *could* occur were to drain attention away from actual catastrophes that *are* occurring.

It is also appropriate, at this juncture, to reflect for a moment on the biggest cause of death and disability of all, namely ageing, which accounts for perhaps two-thirds of the 57 million deaths that occur each year, along with

an enormous loss of health and human capital.¹² If ageing were not certain but merely probable, it would immediately shoot to the top of any list of global catastrophic risks. Yet the fact that ageing is not just a possible cause of future death, but a certain cause of present death, should not trick us into trivializing the matter. To the extent that we have a realistic prospect of mitigating the problem – for example, by disseminating information about healthier lifestyles or by investing more heavily in biogerontological research – we may be able to save a much larger expected numbers of lives (or quality-adjusted life-years) by making partial progress on this problem than by completely eliminating some of the global catastrophic risk discussed in this volume.

Other global catastrophic risks which are either already substantial or expected to become substantial within a decade or so include the risks from nuclear war, biotechnology (misused for terrorism or perhaps war), social/economic disruption or collapse scenarios, and maybe nuclear terrorism. Over a somewhat longer time frame, the risks from molecular manufacturing, artificial intelligence, and totalitarianism may rise in prominence, and each of these latter ones is also potentially existential.

That a particular risk is larger than another does not imply that more resources ought to be devoted to its mitigation. Some risks we might not be able to do anything about. For other risks, the available means of mitigation might be too expensive or too dangerous. Even a small risk can deserve to be tackled as a priority if the solution is sufficiently cheap and easy to implement – one example being the anthropogenic depletion of the ozone layer, a problem now well on its way to being solved. Nevertheless, as a rule of thumb it makes sense to devote most of our attention to the risks that are largest and/or most urgent. A wise person will not spend time installing a burglar alarm when the house is on fire.

Going forward, we need continuing studies of individual risks, particularly of potentially big but still relatively poorly understood risks, such as those from biotechnology, molecular manufacturing, artificial intelligence, and systemic risks (of which totalitarianism is but one instance). We also need studies to identify and evaluate possible mitigation strategies. For some risks and ongoing disasters, cost-effective countermeasures are already known; in these cases, what is needed is leadership to ensure implementation of the appropriate programmes. In addition, there is a need for studies to clarify methodological problems arising in the study of global catastrophic risks.

¹² In mortality statistics, deaths are usually classified according to their more proximate causes (cancer, suicide, etc.). But we can estimate how many deaths are due to ageing by comparing the age-specific mortality in different age groups. The reason why an average 80-year-old is more likely to die within the next year than an average 20-year-old is that senescence has made the former more susceptible to a wide range of specific risk factors. The surplus mortality in older cohorts can therefore be attributed to the negative effects of ageing.

The fruitfulness of further work on global catastrophic risk will, we believe, be enhanced if it gives consideration to the following suggestions:

- In the study of individual risks, focus more on producing actionable information such as early-warning signs, metrics for measuring progress towards risk reduction, and quantitative models for risk assessment.
- Develop and implement better methodologies and institutions for information aggregation and probabilistic forecasting, such as prediction markets.
- Put more effort into developing and evaluating possible mitigation strategies, both because of the direct utility of such research and because a concern with the policy instruments with which a risk can be influenced is likely to enrich our theoretical understanding of the nature of the risk.
- Devote special attention to existential risks and the unique methodological problems they pose.
- Build a stronger interdisciplinary and international risk community, including not only experts from many parts of academia but also professionals and policymakers responsible for implementing risk reduction strategies, in order to break out of disciplinary silos and to reduce the gap between theory and practice.
- Foster a critical discourse aimed at addressing questions of prioritization in a more reflective and analytical manner than is currently done; and consider global catastrophic risks and their mitigation within a broader context of challenges and opportunities for safeguarding and improving the human condition.

Our hopes for this book will have been realized if it adds a brick to the foundation of a way of thinking that enables humanity to approach the global problems of the present era with greater maturity, responsibility, and effectiveness.

