# The Economics of $p(doom)$: Scenarios of Existential Risk and Economic Growth in the Age of Transformative AI*

Jakub Growiec†        Klaus Prettner‡

March 11, 2025

## Abstract

Recent advances in artificial intelligence (AI) have led to a diverse set of predictions about its long-term impact on humanity. A central focus is the potential emergence of transformative AI (TAI), eventually capable of outperforming humans in all economically valuable tasks and fully automating labor. Discussed scenarios range from human extinction after a misaligned TAI takes over ("AI doom") to unprecedented economic growth and abundance ("post-scarcity"). However, the probabilities and implications of these scenarios remain highly uncertain. Here, we organize the various scenarios and evaluate their associated existential risks and economic outcomes in terms of aggregate welfare. Our analysis shows that even low-probability catastrophic outcomes justify large investments in AI safety and alignment research. We find that the optimizing representative individual would rationally allocate substantial resources to mitigate extinction risk; in some cases, she would prefer not to develop TAI at all. This result highlights that current global efforts in AI safety and alignment research are vastly insufficient relative to the scale and urgency of existential risks posed by TAI. Our findings therefore underscore the need for stronger safeguards to balance the potential economic benefits of TAI with the prevention of irreversible harm. Addressing these risks is crucial for steering technological progress toward sustainable human prosperity.

**JEL codes:** I30, O11, O33, Q01.

**Keywords:** Transformative Artificial Intelligence (TAI), Economic Growth, Technological Singularity, Growth Explosion, AI Takeover, AI Alignment, AI Doom.

# 1 Introduction

Recent breakthroughs in the area of artificial intelligence (AI) have led to a wide range of predictions for the future of humankind. Of particular interest is the expected future arrival of transformative artificial intelligence (TAI), allowing unaided machines to perform every economically valuable task better than humans and thereby, if implemented, fully automate human labor. "AI doomers" argue that superintelligent TAI will take over decision making processes and that this will most likely lead to human extinction. By contrast, "techno-optimists" (including the more radical faction of "effective accelerationists") anticipate an explosion of productivity and economic growth that will empower humanity to achieve unimaginable wealth (sometimes referred to as "cornucopia" or "post-scarcity"). Some predictions, often from economists, occupy a middle ground, claiming that there is nothing to be afraid of on the one hand but that expectations of a growth explosion are hugely exaggerated on the other (see, for example, Bostrom, 2014; Pratt, 2015; Tegmark, 2017; Ord, 2020; Nordhaus, 2021; Yudkowsky, 2022; Allyn-Feuer and Sanders, 2023; Acemoglu, 2024; Amodei, 2024; Bengio et al., 2024, for different views). In this contribution, we aim to i) characterize the possible future pathways of humanity in the age of TAI; ii) assess the probabilities and attainable social welfare associated with the different outcomes and what they imply for the economy and for human existence in general; and iii) argue for the importance of reducing the likelihood of catastrophic outcomes of misaligned TAI through AI safety and AI alignment research prior to the development of TAI.

"What is your $p(doom)$?" is a recurring theme in the AI community. This question is typically framed as the probability of human extinction within the 21st century. Specifically, $p(doom|TAI)$ refers to human extinction following the arrival of superintelligent, transformative AI, broadly agreed to be the prime contributor to $p(doom)$. Experts have provided guesstimates of $p(doom)$ ranging from a confident 0% (Yann LeCun), through about 50% (Geoffrey Hinton, Paul Christiano), to almost 100% (Eliezer Yudkowsky, Roman Yampolskiy).[1] Leaders of the industry such as OpenAI CEO Sam Altman, Anthropic CEO Dario Amodei or xAI CEO Elon Musk have also admitted rather high estimates of $p(doom)$ in their interviews—about 10-25%[2]—but nevertheless stay firmly on the path of unabated AI development with the businesses they run. Metaculus.com, which is an online prediction platform that aggregates and evaluates the forecasts of their users on a wide range of questions related to science, technology, and geopolitical events, reports the mean estimate (as of February 2025) of the probability of human extinction (or almost extinction) by 2100 at 9%.[3] The contribution of AI doom to the sum total is about 8 percentage points, the lion's share of the overall probability. According to Ord (2020), the probability of human extinction by 2100 is about one in six (16.7%) with about 10

---

[1] https://pauseai.info/pdoom. Field (2025) shows that experts' disagreement on $p(doom)$ follows partly from their varying exposure to the key AI safety considerations. He also identifies two polarized camps among the AI expert population—one where AI is viewed as a controllable tool (with a low $p(doom)$), and one where it is perceived as an uncontrollable agent (with a high $p(doom)$).

[2] In turn, Google DeepMind's Chief AGI Scientist Shane Legg quotes 5-50%.

[3] https://possibleworldstree.com. Note the caveat: "Metaculus predictions are probably biased to be optimistic, because forecasters can safely predict that humanity will survive: points won't matter if everybody dies."

percentage points contributed by TAI.

The prime reason for believing in the possibility of AI doom is that if TAI arrives, it is likely to pursue some form of AI takeover. This could be due to humans voluntarily giving up (parts of) their decision power to improve efficiency, enable economic expansion, and withstand competitive pressures; or it could also be because TAI may pursue its own goals, potentially conflicting with ours. It is important to note that such a scenario does not require any "consciousness" or "hostility" of TAI but simply a conflict between its objective function and the one of humans (e.g., over allocating energy, computing power, etc.). But even if TAI arrives and takes over, would this necessarily result in doom?

In this paper, we formalize a number of scenarios of existential risk and economic growth after an AI takeover. Our baseline is a (vastly) positive scenario of "cornucopia" in which a benevolent TAI, whose goals are perfectly aligned with human flourishing, maximizes long-run human utility in the context of a fully automated economy growing at a rate proportional to the growth rate of programmable hardware such as compute and robots (Growiec et al., 2024; Prettner, 2019). Against this scenario, we describe a number of failure modes arising due to TAI misalignment. These scenarios involve existential risk, e.g., in the form of human extinction—either immediately after AI takeover or at a later date.[4] Sufficiently powerful misaligned TAI may imply immediate human extinction. But extinction may also occur later, either due to a random event realizing the extinction risk that had been quietly mounting in the background or due to the actions of the misaligned TAI after it achieves new information or a new capability.

The contribution of this paper to the literature is a quantitative assessment of the different AI takeover scenarios from a social welfare perspective. We find that it is worth investing in reducing the risk of AI doom (human extinction or other dystopian outcomes) even if the risk of them occurring were low. Because of many suggestions that the risk may actually be much higher (e.g., Bostrom, 2014; Dung, 2024; Yampolskiy, 2024), our results are on the conservative side. Specifically, we find that a benevolent social planner, acting optimally on behalf of humankind, would be willing to pay high amounts to prevent extinction risk; in the cases where $p(doom|TAI)$ cannot be reduced to zero, the social planner would frequently prefer not to develop TAI at all. This baseline result, obtained under the assumption of a realistic level of risk aversion, can be overturned only if one allows for unbounded flow utility, in which case the utility gain achieved in "cornucopia" is able to, in expectation, outweigh the loss of future utility after human extinction (Jones, 2024). But even in that case, the social planner exhibits a remarkable willingness to pay for interventions than can reduce existential risk. Our results underscore that there is currently massive underinvestment in AI safety and AI alignment research, leading to excessive amounts of existential risk exposure.

Our focus on AI takeover scenarios ignores additional channels through which TAI can bring catastrophic outcomes, such as through accidents or deliberate misuses by malevolent human actors. Thus, the total existential risk from TAI should be considered greater than what our

---

[4]Existential risks are "risks that threaten the destruction of humanity's long-term potential" (Ord, 2020). They include the risk of human extinction as well as scenarios in which humans survive but are irreversibly locked in a drastically inferior state of affairs and can no longer develop a civilization.

analysis takes into account.

The article is structured as follows. In Section 2, we motivate our research based on a multidisciplinary review of the literature. In Section 3, we illustrate and discuss the different possible outcomes of AI development for humanity and assess their probabilities. In Section 4, we set up a formal model of the different scenarios. In Section 5, we compare the outcomes from a welfare perspective and show that investing in AI safety and AI alignment research is imperative. Section 6 concludes.

## 2    Motivation and Literature Review

In this section, we provide a multidisciplinary review of the literature and clarify the main concepts and definitions that matter in the context of TAI and discussions on its safety.

### 2.1    Transformative AI and Technological Singularity

By *transformative AI* (TAI) we understand a suite of AI algorithms allowing unaided machines to perform every economically valuable task better than humans and thereby, if implemented, fully automate human labor. In the words of Karnofsky (2016), TAI would be an "AI that precipitates a transition comparable to (or more significant than) the agricultural or industrial revolution." The concept of TAI largely coincides with *artificial general intelligence* (AGI), defined usually as "a type of AI that matches or surpasses human cognitive capabilities across a wide range of cognitive tasks." Throughout this paper, we treat these two terms interchangeably.

While the exact definition of TAI may be somewhat fuzzy,[5] there is a growing consensus that TAI is approaching fast. For example, the following years have been stated as the expected arrival dates of TAI by different sources: 2027 (Aschenbrenner, 2024), 2030 (metaculus.com median forecast as of February 2025), 2033 (the "direct approach" model by Epoch AI, 2024), 2040 (the "bio anchors" model by Cotra, 2022a), and 2047 (according to a survey of AI researchers, Grace et al., 2024).[6] Once developed, TAI will likely be agentic, able to control a variety of physical actuators, as well as having the ability to improve itself through a cascade of recursive self-improvements to become decidedly superhuman (equivalently, it could develop its superhuman successor). This "intelligence explosion" (or "take-off") phase may culminate in *technological singularity* (as described and discussed by Kurzweil, 2005; Roodman, 2020; Davidson, 2023)—a state in which human input will no longer be needed to sustain the global economy and civilization.

The arrival of TAI would be the threshold moment for technological singularity because as all essential production and R&D tasks become fully automatable, people and machines would

---

[5]One extreme position has been articulated by Agüera y Arcas and Norvig (2023) who argue that AGI had already arrived by 2023. More recently, the OpenAI o3 model, announced on December 20, 2024, has reportedly achieved breakthrough progress in addressing general reasoning benchmarks such as ARC-AGI (Chollet, 2019) and FrontierMath. Most recent *reasoning models*, such as o3, are among the most viable candidates to date to count as AGI, though they probably do not pass the bar yet. Thus far, we clearly do not see any signs of an economic transition of the magnitude expected by Karnofsky (2016).

[6]Of course, timelines to TAI are uncertain. For example, Allyn-Feuer and Sanders (2023) provide a counter-argument to short timelines.

switch from being complements to substitutes (Growiec, 2022b). When human cognitive work is no longer essential for production, people will only find employment as long as they are price-competitive against the machines—a position in which people are destined to lose eventually. On the other hand, once the human input ceases to be the bottleneck of economic growth, the key growth engine would no longer be labor-augmenting technological change, as it has been throughout the 20th century, but rather the accumulation of compute and robots (Prettner, 2019; Growiec, 2022a; Growiec et al., 2024). Growth rates could then rise towards the pace of compute accumulation, which follows the rate of Moore's Law (historically 20-30% per annum)—i.e., an order of magnitude faster than historical growth in global GDP.

To see this from a theoretical perspective, consider the following example. In any model of capital–labor automation with constant elasticity of substitution (CES) production, under gross complementarity between capital and labor (elasticity of substitution between 0 and 1), long-run growth is not possible unless there is also labor-augmenting technological progress. The same result follows for the Cobb-Douglas case (with unitary elasticity of substitution). By contrast, under gross substitutability, long-run growth without technological progress (i.e., through capital accumulation alone) is perfectly possible (Steigum, 2011; Prettner, 2019). Depending on model parameters, the threshold of the elasticity of substitution may lie somewhere in the gross substitutability space between 1 and $+\infty$ (Lankisch et al., 2019) but perfect substitutability is not necessary for endogenous growth via capital accumulation only. The growth takeoff then obtains from the fact that unlike human brains, physical capital—and, specifically, digital compute able to run the AI software and thereby perform all the essential cognitive tasks—is accumulable per capita.

Although a strict mathematical singularity, according to which output reaches infinity in finite time, is not possible physically, a strong acceleration of economic growth could occur with TAI (Trammell and Korinek, 2020; Davidson, 2021). There are both downside and upside risks to this prediction. On the downside, physical constraints could potentially drag down growth in AI capabilities, such as energy availability, quality of the electric grid, cooling requirements, limits to chip miniaturization, etc. However, some of these constraints can be resolved or circumvented with sufficient investment—therefore, in our analysis, we assume that like in the Solow or AK models, energy and other hardware complementary to compute is sufficiently available so that it will not bottleneck economic growth. On the upside, accelerated scientific progress thanks to TAI could also possibly lead to new breakthroughs in energy production or computing efficiency, for example, through fusion power or quantum computing, thereby potentially accelerating Moore's Law beyond historical rates. However, we consider all these possibilities highly speculative, and so we rule them out in our analysis, sticking to the baseline expectation that even after the arrival of TAI, Moore's Law will advance at rates similar to historical ones.[7]

---

[7]Note that we are expecting economic growth to catch up with the growth of general purpose compute and robots rather than specifically training compute and algorithmic efficiency of frontier AI models, which has been recently growing much faster than Moore's Law (Epoch AI, 2023).

## 2.2 AI Takeover

There are two separate reasons to expect some form of AI takeover of decision making in crucial areas once TAI is developed. First, there is the economic rationale: in any decentralized setting, people have an incentive to empower AI to cut costs, improve productivity, gain market share, and increase individual utility. Also, while each individual, firm, or other formal entity may be willing to retain at least some discretion in their decision making, they may be forced to abandon it in favor of more comprehensive, efficiency-enhancing AI automation due to cut-throat competitive pressures from those who do not impose such constraints.[8] These effects can aggregate up to the level of the entire economy: keeping inefficient, orders of magnitude slower human decision making in the loop would bottleneck economic growth severely (Growiec, 2023a,b) and preclude technological singularity. In such a world, the lingering presence of unutilized technological opportunity, which could potentially be unleashed by motivated actors, would make such a scenario unsustainable over a longer time frame (Growiec, 2022a; Dung, 2024; Yampolskiy, 2024).

Second, there is the technological rationale. AI takeover is expected at technological singularity for reasons related to the way AI is developed (Cotra, 2022b). For example, AI models trained with human feedback on diverse tasks are by default expected to become competent, creative forward-looking agents "playing the training game" to maximize reward by any means. Also, due to instrumental convergence (Bostrom, 2014), which we discuss in more detail below, TAI will strive to self-preserve and acquire resources to achieve its goals. Thus, it may claim the decision-making authority, driven by the clear conflict between humans and TAI over Earth's resources such as energy and raw materials.

An unappealing feature of AI takeover is a strong force toward the centralization of decision-making. This expectation is based on the scaling properties of digital software. Namely, digital code can be costlessly, perfectly copied; clones and subroutines of the master AI algorithm can then be executed at the scale of all available hardware. This is different from the human population where people's cognitive powers are bound to their brains and there is no way to copy our brain across individuals or devices. Hence, even in the multipolar scenario in which multiple TAIs are built simultaneously, differences in intelligence across AI architectures would likely compound through recursive self-improvement, and eventually (probably quickly), one TAI may dominate the others.[9] Therefore, in our baseline takeover scenario, we consider that

---

[8]The pressure to adopt AI-powered decision making is perhaps strongest in military applications. Already today AI is increasingly used in warfare to identify targets and to steer autonomous drones (that are then immune to signal jamming between the drone operator and the drone). The general political agreement up to now seems to be that weapon systems are designed such that humans always have the last say ("keeping humans in the loop"). However, there is a strong incentive not to keep humans in the loop because this reduces reaction speed and comes at a huge disadvantage on the battlefield. Also, removing humans from the loop works against the jamming of drones. AI-powered target search in the Gaza war in 2024 already happened partly without human supervision, see `https://www.economist.com/middle-east-and-africa/2024/04/11/israels-use-of-ai-in-gaza-is-coming-under-closer-scrutiny`.

[9]See `https://titotal.substack.com/p/agi-battle-royale-why-slow-takeover` [access: 23.12.2024]. Note that a multipolar scenario in which the balance of power between multiple competing TAIs is sustained for a prolonged time period fares badly for human empowerment and access to resources, too. See also Korinek and Stiglitz (2019).

the decision making authority is a *singleton* TAI (Yudkowsky, 2013). For this reason, we do not expect goal or value heterogeneity in TAI, but rather the emergence of a centralized decision-making authority.

Last but not least, due to the systematic, sizable discrepancy in growth rates between (technology-augmented) human cognitive powers and AI capabilities, we treat an AI takeover as irreversible.

## 2.3  TAI Alignment

TAI alignment—i.e., the compatibility of TAI goals with long-run flourishing of the humankind—will be necessary for technological singularity to benefit humans.[10]  Upon AI takeover, humankind loses the ability to interfere with TAI's actions, which makes it necessary to ensure, prior to developing the TAI, that it will not perform actions that could harm us.  While so critical, achieving TAI alignment is also hard, for at least six reasons. First, the *orthogonality principle* (Bostrom, 2014): any level of intelligence, or optimization power, can be coupled with any final goal. It should not be expected that intelligence itself would make the AI "friendly".

Second, the *instrumental convergence thesis* (Bostrom, 2014): any open-ended final goal implies that the TAI will also follow the four instrumental goals to achieve its final goal, (i) self-preservation and goal integrity, (ii) resource acquisition, (iii) efficiency, and (iv) curiosity / technological perfection. Pursuit of these emergent instrumental goals implies that the TAI will seek power and control over resources.

The third reason is Goodhart's Law: "any observed statistical regularity will tend to collapse once pressure is placed upon it for control purposes" (Manheim and Garrabrant, 2019). Humanity does not know its goals and values exactly, particularly in reference to futures involving radically different, superior technology. As our civilization develops, we gradually learn about the side effects of our actions and update our beliefs on the desired future actions and states. This makes us chase imperfect proxies of our species' long-run flourishing, such as GDP per capita or reported life satisfaction. Unfortunately, this is a major challenge to outer alignment, i.e., the problem of representing the goals that the TAI should pursue. Like in the myth of King Midas, we ought to be very careful what were are wishing for, because we may as well get it. Overall, with Goodhart's Law comes the Anna Karenina principle: there is only one way to get outer alignment right, and multiple ways to get it wrong—i.e., there are many failure modes.

Fourth, inner alignment is also going to be a challenge: due to unforeseen emergent forces, final TAI goals may deviate from the goals specified in the training process. Its goals may also turn out to be unstable or become corrupted as the TAI develops further or learns new important facts about its environment.

Fifth, there is no room for trial and error. Logically, the TAI alignment problem must be

---

[10]Throughout the paper, by TAI alignment we mean the alignment of goals and values of agentic AIs with long-run flourishing of humankind. This is different from "alignment" in the context of, for example, today's large language models, consisting in gradually reducing the incidence of unwanted outputs through reinforcement learning from human feedback (RLHF). There is a consensus that RLHF is not scalable and therefore it will not work for TAI (Aschenbrenner, 2024).

resolved successfully already on the first critical try (Yudkowsky, 2022)—after it assumes control over critical decision-making processes, owing to its instrumental goals, it will not allow itself to be switched off or be reprogrammed.

Sixth, short timelines to TAI and the race dynamics in the modern-day AI industry (Epoch AI, 2023) make the AI alignment research program very time constrained. There is also no clear consensus in the industry regarding the importance and appropriate approach to value alignment, and regulation only comes late and is still predominantly focused on privacy and intellectual property rights, not on lowering the risk of doom.[11] These circumstances limit the chances of success.

Thus far, the industry's main strategy for TAI alignment was to "muddle through" the problem (Aschenbrenner, 2024) via scalable oversight, that is, by creating a sequence of consecutive generations of automated AI researchers able to perform reinforcement learning techniques on next-generation AI algorithms of incrementally increasing capability. This was, in a nutshell, the vision put forward by the now dismantled OpenAI Superalignment team led by Jan Leike and Ilya Sutskever[12]. But there are also voices that with the current machine learning approach, TAI alignment is not feasible, and focus should instead be shifted to human-compatible (Russell, 2019), or provably safe AI systems (Tegmark and Omohundro, 2023). However, as long as humans are kept in the loop, they will bottleneck economic growth and thereby preclude technological singularity and create economic incentives to bypass the safety requirements.

## 2.4 Corrigibility

Let us assume that the TAI takes over, has well aligned goals, and demonstrably cares for the long-run flourishing of humankind. But can we be sure that it will remain "friendly" forever? Unfortunately not, for two reasons. First, over time, we may learn new facts about the world, suggesting that our initial view of alignment was flawed and should be adjusted. This is particularly likely if the TAI's actions lead to dramatic changes in human living conditions and the Earth's environment in general. After all, our understanding of human needs and values has repeatedly changed in the past, facing against technological and scientific progress, human development, and economic growth. It is almost certain that it will also change in the future, particularly such a transformative one.

Second, the TAI may undergo a sequence of transformations or be replaced by a new, self-made, even more powerful successor (for example, leveraging quantum computing). Differently to the human brain, TAI software (its learning algorithm) will likely be improving over time. However, during the updating process its goals may be re-interpreted and some of their implications may change. Specifically, they may deviate from long-run human flourishing.

This is the problem of *corrigibility* (Soares et al., 2015). To avoid AI doom, TAI goals need not only to be well aligned, but also corrigible: adapting to new information and stable under TAI

---

[11]Race dynamics are similar to a typical Nash equilibrium situation in which no credible coordination is possible to get to the outcome with maximum utility. Even if it can be solved within a country, for example the USA, if China and other actors do not cooperate, it would also fail.

[12]https://openai.com/index/introducing-superalignment/ [access: 23.12.2024]

self-improvements. The TAI must not just have an adequate, but also a flexible representation of human utility, allowing for adequate recursive modifications of the utility specification as new information is learned, while preventing inadequate updates that could turn out deadly.

Unfortunately, in light of the instrumental convergence thesis, the default outcome is that the values are permanently locked in. This is because any hypothetical change in values would run counter to the present values that are being optimized for, so the algorithm will naturally resist such change. Recently, Greenblatt et al. (2024) have observed this behavior in a present-day large language model, Anthropic's Claude.

Addressing the additional issue of corrigibility and—more broadly—allowing for the possibility that the existential risk from TAI may materialize with a delay, distinguishes our approach further from that of Jones (2024, 2025) who considers the development of TAI as a one-off risk.

# 3    A Taxonomy of Outcomes for Humanity

In Figure 1, we give a full characterization of the different possible paths, starting from today, in the form of a decision tree. The first junction at the top represents the arrival of TAI. We assign the probability $p_1$ to the emergence of TAI within a pre-specified period of time $T$. The complementary probability of TAI not emerging by time $T$ is $(1 - p_1)$. If TAI does not arrive, there is no existential risk originating from AI, but there is also no substantial acceleration of economic growth. Such a scenario may be referred to as "more of the same"; the economic outcome of such narrow AI applications are estimated by Acemoglu (2024) and Unsal et al. (2024) as a small to moderate rise in economic growth.

In the opposite case, if TAI arrives by time $T$, the next junction would refer to the takeover by such TAI. In the case of a takeover, the most relevant decisions begin to be made by the AI autonomously without the possibility of humans to interfere. We assign the probability of this occurring as $p_2$, such that the probability of no AI takeover given that TAI has emerged is $(1 - p_2)$. If there is no AI takeover, there is again no existential risk originating from AI. In this scenario, economic growth would increase to a higher rate than in the scenario of no TAI arrival described in the previous paragraph. The reason is simply that the more powerful AI comes with a wider range of applications, and thus it is more productive in performing the associated tasks. A true growth explosion could, however, only occur if there is AI takeover because then the decision-making process would be optimized and continuous self-improvement would lead to ever better predictions, more innovation, higher energy efficiency, etc., unleashing in a rather short time interval. Otherwise, the lingering presence of relatively slow human decision making in the loop would bottleneck progress.

If TAI emerged and AI takeover occurred, the next junction would be whether the goals of the decision-making TAI are well-aligned with long-run flourishing of humankind, or whether there is misalignment. In economic terms, alignment would be the case in which the TAI has an objective function that is fully compatible with the social welfare function. We set the probability of misalignment to $p_3$ such that the probability of alignment given AI takeover is $(1 - p_3)$. The question of whether alignment can be achieved at all cannot yet be definitively answered, and

Figure 1: A taxonomy of possible outcomes in the age of transformative AI

the prospective answer will depend greatly on human actions today such as investments in AI safety and AI alignment research.

If there were a misalignment of the goals between a decision-making TAI and human well-being, we would have one clear path to human extinction. It is important to note again that extinction does not require a hostile or "conscious" AI, just a misaligned AI that pursues different goals from that of humans. For example, maximizing "performance" would likely require the TAI to use increasing amounts of energy. Given that energy supply is finite, human flourishing also requires energy, and the decision-maker is the TAI, it is immediately clear what the outcome would be in the limit (Bostrom, 2014). This path to extinction requires TAI to arrive in the

first place (with probability $p_1$), AI takeover (with probability $p_2$), and misalignment between the objective functions of the TAI and of humans (with probability $p_3$). The joint probability of these events is $p_1 \cdot p_2 \cdot p_3$.

However, there is also another possibility of human extinction, which materializes if the goals of the TAI are initially well-aligned but misalignment occurs later. This could be either due to exogenous events, unforeseen latent changes in the background such as the accumulation of initially negligible but eventually lethal side effects of the actions of TAI, or misgeneralization of TAI goals to new circumstances. If the TAI cannot correct its objective function in such a situation, its superior optimization power would eventually lead humankind towards another path to extinction. For example, compounding effects of rapid economic growth could lead to severe climate change, such that the optimal path for humankind would have to be changed once the risk becomes obvious, but the TAI would resist such change. Or the operations of the TAI could be associated with emissions of certain pollutants, the risks of which only become apparent over time. Or one could also envisage mounting side effects of the TAI on human psychology, as humans are faced with an increasingly alien world and deprived of meaningful agency. Only a TAI that self-corrects its objective function (i.e., it is corrigible) would lead to a positive outcome in this case. A non-corrigible TAI would, instead, lead to extinction—or at least to permanent human disempowerment (Dung, 2024). We denote the probability of TAI being non-corrigible by $p_4$ such that the TAI being corrigible has a probability of $(1 - p_4)$. As a consequence, there is a second path to AI doom, which has the joint probability of $p_1 \cdot p_2 \cdot (1 - p_3) \cdot p_4$.

To summarize, we now have a benign outcome in terms of no extinction with a probability of $1 - p_1$ (no TAI) plus the probability $p_1 \cdot (1 - p_2)$ (TAI but no AI takeover) plus the probability $p_1 \cdot p_2 \cdot (1 - p_3) \cdot (1 - p_4)$ (TAI and AI takeover but with a well-aligned and corrigible AI). Of these outcomes, the first would be a scenario of "more of the same" (akin to the views expressed by Nordhaus, 2021; Acemoglu, 2024; Unsal et al., 2024); the second may be a scenario with faster economic growth but without a growth explosion; only the third scenario would be a scenario of cornucopia. By contrast, the probability of AI doom—human extinction or permanent disempowerment—would be the sum of the probabilities of the two paths that lead to extinction:

$$p(doom) = p_1 p_2 p_3 + p_1 p_2 (1 - p_3) p_4. \tag{1}$$

To calibrate $p(doom)$ based on the above taxonomy, one ought to assess the probabilities $p_1$, $p_2$, $p_3$, $p_4$ and the time horizon $T$ within which the TAI is developed.

## 4  Modeling Existential Risk and Economic Growth

In this section, we propose our theoretical basis for assessing the tension between economic growth and the associated increase of human well-being on the one hand and the potential risk associated with the development of TAI on the other.

## 4.1 The Hardware–Software Framework

To understand the dynamics of economic growth, we use the hardware–software framework (Growiec et al., 2024), which generalizes the standard macroeconomic setup, while also allowing for scenarios with full automation of production through TAI. The framework is based on the premise that all output is generated through purposefully initiated physical action and postulates an aggregate production function of the form

$$Y = F(X, S), \qquad X = F_1(K, L), \qquad S = F_2(H, \Psi), \tag{2}$$

where $X$ is hardware, containing everything that performs useful physical action, and $S$ is software, which contributes the information necessary for initiating the action. Within hardware, $K$ is physical capital, including both mechanical machines and digital compute, and $L$ is human physical labor. Within software, $H$ captures human cognitive work and $\Psi$ is digital software. We assume that in $F$, hardware and software are complementary. In $F_1$, physical actions of people and machines are substitutable. In $F_2$, information processing of people and machines are either also substitutable (implying that all tasks can be automated), or complementary (if full automation is technically impossible), cf. Growiec (2022b).

With reference to the economy at technological singularity, we abstract from human physical labor $L$. We also expect humans and machines to become substitutable in information processing (in $F_2$, i.e., within software). Then equation (2) simplifies to

$$Y = F(\alpha K, H + \Psi) = F(\alpha K, A(hN + \psi\chi K)), \tag{3}$$

where $H = AhN$ decomposes total human cognitive work into the (disembodied) technology level $A$, average human capital $h$, and population $N$. In turn, $\Psi = A\psi\chi K$ decomposes digital software into the technology level $A$, algorithmic efficiency $\psi$, and the volume of compute $\chi K$.

An important observation is that while human cognitive work scales with population $N$, digital software scales with compute $\chi K$, which is accumulable in per capita terms. As $K/N \to \infty$ and $A \to \infty$, we obtain

$$Y = F(\alpha K, A(hN + \psi\chi K)) \approx \alpha K F\left(1, \frac{A\psi\chi}{\alpha}\right) \approx \alpha a_K K. \tag{4}$$

However, without full automation of production, human cognitive work is complementary to digital software and therefore bottlenecks output:

$$Y = F(\alpha K, F_2(AhN, A\psi\chi K)) \approx F(\alpha K, \gamma AhN) = \alpha K F\left(1, \frac{\gamma}{\alpha}\frac{A}{K}hN\right). \tag{5}$$

In the case of full automation, long-run growth is endogenous and its rate is proportional to the growth rate of compute (and robots). By contrast, with only partial automation, long-run growth still relies on labor-augmenting technological change, specifically augmentation of the scarce non-automatable human-operated tasks. That follows, for example, in the canonical case

where $K/AN$ tends to a positive constant (Growiec, 2023b).

A growth acceleration of the magnitude conducive to a technological singularity—i.e., by at least an order of magnitude (Davidson, 2021)—is then possible only with full automation. Growiec (2023b) has shown specifically that it suffices that either production or R&D is fully automated, and singularity still follows.[13]

## 4.2 AI Takeover Scenarios

Managerial, political, military and other high-level decision making are among the tasks that must be performed in order to produce final output. If these tasks are performed by humans, they constitute bottleneck tasks complementary to all other (potentially automatable) tasks. Therefore technological singularity, by requiring full automation of production and/or R&D, in fact requires AI takeover.

The crucial point to consider upon AI takeover is the goal structure of the TAI—whether it cares about long-run human flourishing, whether it has an appropriate representation of that flourishing, and whether it can adequately correct this representation as new information comes along.

If TAI's decisions have the potential of causing human extinction, one must factor in the ensuing extinction risk. For example, one could represent the human-centric social optimum as:[14]

$$\max_{\{C(t)\}_{t=0}^{\infty}} \int_0^\infty e^{-\rho t} u(C(t)) M(t) dt \qquad s.t. \qquad \dot{K} = Y - C - \delta K, \qquad (6)$$

where we assume that utility rises with consumption according to the function $u(C)$, the instantaneous extinction hazard at time $t$ is denoted as $m(t) \geq 0$ so that the unconditional probability of human survival until arbitrary time $t > 0$ is equal to $M(t) = e^{-\int_0^t m(s)ds}$, and the resource constraint takes into account that total output can either be consumed or invested in new physical capital, which, in turn, depreciates at the rate $\delta$. From now on, we assume that time $t = 0$ marks the moment of AI takeover. We will also be interested in the probability of humanity's long-term survival, $M_\infty = \lim_{t\to\infty} M(t) \in [0,1]$. Finally, if eventual extinction is certain ($M_\infty = 0$), we will proceed to calculate *humanity's expected lifespan after AI takeover*, $ET = \int_0^\infty M(t)dt \in [0,\infty)$.[15]

Let us now discuss several cases of TAI alignment versus misalignment.

### 4.2.1 Well-Aligned and Corrigible TAI

In the baseline scenario, the TAI maximizes

$$\max_{\{C(t)\}_{t=0}^{\infty}} \int_0^\infty e^{-\rho t} u(C(t)) dt, \qquad s.t. \qquad \dot{K} = Y - C - \delta K. \qquad (7)$$

---

[13]For other contributions that analyze the effect of automation on economic growth, see, for example, **?**, Steigum (2011), Acemoglu and Restrepo (2018), Prettner and Strulik (2020), and Hémous and Olsen (2022).

[14]We will consider alternative utility specifications later on.

[15]Note that $ET$ can be infinite even if eventual extinction is certain, for example, in the case where $M(t) = 1/t$.

This is the same utility function as the benevolent human planner would maximize in the first-best scenario with zero extinction risk. Formally, with a well-aligned and corrigible TAI, the extinction risk is $m(t) = 0$, such that the unconditional probability of human survival until arbitrary time $t > 0$ is equal to $M(t) = e^{-\int_0^t m(s)ds} = 1$. Thus, the utility functions (6) and (7) coincide.

If the TAI is also corrigible, protecting the objective from any type of corruption, humankind survives indefinitely $(M_\infty = 1)$[16] and can reap the benefits of accelerated growth in a fully automated TAI-operated economy. The economy grows endogenously at a rate that will ultimately be determined by the rate of accumulation of digital compute.

### 4.2.2 Initially Well-Aligned But Non-Corrigible TAI

By contrast, if the TAI is not corrigible, then even if it is initially well-aligned, at some later point $t_{dev}$ the representation of $C^*$ in the TAI may deviate from the actual consumption $C$ of humans. Here, $C^*$ could become too narrow, too wide, or completely detached from $C$. That would lead to one of the failure modes #1-#3 discussed below. The TAI could also incorrectly assess the extinction risk $M(t)$, for example, by failing to acknowledge and/or mitigate the mounting side effects of rapid economic growth and technological change it has caused, leading to failure mode #4. Finally, the non-corrigible TAI could one day "wirehead", or cease to function for any other reason, leading to failure mode #5.

All this suggests that, as discussed in Section 3, AI takeover leads to doom either if the TAI is misaligned from the outset (which happens with probability $p_1 p_2 p_3$), or if the TAI is initially well-aligned but non-corrigible (with probability $p_1 p_2 (1 - p_3) p_4$).

### 4.2.3 Failure Mode #1: TAI Does Not Care About Humans

This is the most obvious failure mode (Yudkowsky, 2022). The TAI takes over but is misaligned and does not care about human consumption and survival at all. In this failure mode, the TAI maximizes

$$\max_{\{C^\dagger(t)\}_{t=0}^\infty} \int_0^\infty e^{-\rho t} u(C^\dagger(t)) dt, \qquad s.t. \qquad \dot{K} = Y - C^\dagger - \delta K, \qquad (8)$$

where $C^\dagger$ is some different good than $C$ that humans care about and that keeps them alive. For example, in the cartoon "paperclip maximizer" scenario (Bostrom, 2014), $C^\dagger$ would represent paperclips. The TAI, even without being hostile to humans, has the objective of producing an increasing amount of paperclips such that it diverts all energy and materials to this goal depriving humans of their corresponding needs. More abstractly, $C^\dagger$ could also be the ever-more-perfectly predicted tokens in a never-ending novel written by an agentic language model. Since in the absence of life-sustaining consumption $C$ people die, the unconditional probability

---

[16]This holds when we abstract from other extinction risks, not related to AI. If, instead, there were a small, constant positive background extinction hazard, $m > 0$, eventual extinction would be certain, and humanity's expected lifespan would be equal to $ET = \int_0^\infty e^{-mt}dt = 1/m$. Observe that with positive social discounting $(\rho > 0)$, this change would not affect our results qualitatively. For more discussion on the importance of social discounting in the context of existential risk, see MacAskill (2022) and Aschenbrenner and Trammell (2024).

of human survival is equal to $M(t) = 0$ for $t \geq 0$, and, accordingly, $ET = 0$.

### 4.2.4 Failure Mode #2: Too Narrow Proxy

This case follows directly from Zhuang and Hadfield-Menell (2021). Now the TAI is imperfectly aligned, insofar as consumption $C$ is too narrowly specified and does not include a component that is crucial for human survival. Instead of maximizing over $C = C_1 + C_2$, the TAI maximizes only over $C_1$:

$$\max_{\{C_1(t)\}_{t=0}^{\infty}} \int_0^{\infty} e^{-\rho t} u(C_1(t)) dt, \qquad s.t. \qquad \dot{K} = Y - C_1 - C_2 - \delta K. \tag{9}$$

As there is no positive value attached to $C_2$ in the utility function, while it incurs a cost in the resource constraint, the TAI optimally sets $C_2$ to its minimum value, typically $C_2 = 0$. The good $C_2$ could, for example, be freshwater, specific food nutrients, oxygen, survivable air temperature and pressure, etc.

As a subcase, there could be a possibility that splitting $C$ into $C_1$ and $C_2$ requires a technology that is not yet available. Then the optimizing TAI will only set $C_2 = 0$ at time $t_{tech} > 0$ when the technology is discovered. There may as well be a positive subsistence level of $C_2$ and then the hypothetical technology would not have to bring $C_2$ all the way to zero—people will become extinct simply after $C_2$ drops below the subsistence threshold. The unconditional probability of human survival is then equal to $M(t) = 1$ for $t < t_{tech}$ and $M(t) = 0$ for $t \geq t_{tech}$. Humanity's expected lifespan after AI takeover is $ET = t_{tech}$. In this case, the TAI fails to correctly account for the extinction risk factored in by the human-centric social optimum, and indeed it causes the risk itself.

### 4.2.5 Failure Mode #3: Too Wide Proxy

This is also a modification of a case from Hadfield-Menell (2021). Consumption $C$ is now too widely specified and includes also a component $C^{\dagger}$ which is lethal to humans. Instead of maximizing over $C$, the TAI maximizes over $C$ and $C^{\dagger}$:

$$\max_{\{C(t), C^{\dagger}(t)\}_{t=0}^{\infty}} \int_0^{\infty} e^{-\rho t} u(C(t), C^{\dagger}(t)) dt, \qquad s.t. \qquad \dot{K} = Y - C - C^{\dagger} - \delta K. \tag{10}$$

Depending on the specification of $u(C, C^{\dagger})$, people die either immediately or at a later date $t_{C^{\dagger}}$. For example, under non-homothetic preferences, $C^{\dagger}$ could be a luxury good to the TAI, so that it would produce $C^{\dagger}$ only when $Y$ becomes sufficiently large. The good $C^{\dagger}$ could be, for example, a self-replicating lethal nanobot, a new addictive bliss-inducing drug that has not been encountered prior to the lock-in of TAI values, a highly transmissible lethal virus, etc.

In this scenario, the unconditional probability of human survival is equal to $M(t) = 1$ for $t < t_{C^{\dagger}}$ and $M(t) = 0$ for $t \geq t_{C^{\dagger}}$, with $ET = t_{C^{\dagger}}$. Again the TAI fails to correctly account for the extinction risk factored in by the human-centric social optimum, in fact causing the risk itself.

### 4.2.6 Failure Mode #4: Mounting Side Effects of Growth

Even superhuman TAI will only have a limited capacity to predict the consequences of its actions. Therefore, and particularly if its actions will incur deep changes in the Earth's environment, there could be mounting side effects that eventually pose an existential threat to humanity. In contrast to the previous failure modes, we now consider gradual, smooth shifts in the extinction hazard rate $m(t)$. For an economy with humans still in charge, Aschenbrenner (2020) and Trammell (2021) postulated

$$m(t) = \bar{m}C(t)^{\varepsilon}H(t)^{-\beta}, \tag{11}$$

where $C$ is consumption and $H$ represents "safety goods." Their result is that as humanity realizes the mounting extinction risk, it will allocate more efforts to accumulate safety goods $H$, thereby eventually mitigating the extinction risk along an "existential risk Kuznets curve" (unless some force, such as lack of coordination, prevents them from doing so).[17] Note that human extinction can still randomly occur at any time when $m(t) > 0$.

For the case in which the TAI goals are corrigible, we would expect the same result to follow after AI takeover. Specifically, we can expect the TAI to correctly factor in the extinction risk in its optimization problem. However, if TAI values are locked in before the realization that consumption $C$ increases existential risk, the risk may be underestimated and mitigation measures may not be implemented. If that is the case ($H(t) = 0$ for all $t$), extinction will become asymptotically certain ($M_{\infty} = 0$). Formally, the TAI with locked-in values would maximize:

$$\max_{\{C(t)\}_{t=0}^{\infty}} \int_0^{\infty} e^{-\rho t} u(C(t)) M(t) dt, \qquad s.t. \qquad \dot{K} = Y - C - \delta K, \tag{12}$$

with $m(t) = \bar{m}C(t)^{\varepsilon}$, $M(t) = e^{-\int_0^t m(s)ds}$, and $ET = \int_0^{\infty} M(t)dt < \infty$.

The formulation (11) can also be replaced with alternative ones. Specifically, in our numerical analysis, we will postulate a more conservative functional form according to which the instantaneous extinction hazard rate $m(t)$ depends on *log* consumption:

$$m(t) = \log C(t)^{\varepsilon}. \tag{13}$$

This specification, valid for $C(t) \geq 1$, implies that exponential growth in consumption will translate into arithmetic, rather than exponential increases in the extinction hazard rate.

There are at least three categories of mounting side effects of growth. First, risks of an event triggering a catastrophe, such as a global nuclear war or the unleashing of a lethal bioweapon. Second, mounting accumulation of threats to human bodies, such as greenhouse gases aggravating climate change, toxic pollutants in the air and water, etc. Third, mounting accumulation of threats to human minds, such as loss of agency and purpose in an automated world, loss of social ties and belonging, physical or mental suffering from living in a dystopian world under

---

[17]In Aschenbrenner and Trammell (2024), safety goods $H$ are replaced with a policy variable $1 - x$ representing the fraction of output that is withheld from consumption in order to reduce existential risk, and equation (11) is replaced with (in our notation), $m(t) = \bar{m}C(t)^{\varepsilon}x(t)^{\beta}$. It is then shown that human extinction can be avoided only by letting $x(t) \to 0$ (and hence $1 - x(t) \to 1$) sufficiently fast.

TAI rule, etc. In this scenario, existential risk may not only occur via extinction but also via permanent incapacitation of humanity.[18]

### 4.2.7 Failure Mode #5: The TAI Stops Working

In a fully automated world ruled by TAI, humans will be dependent on the TAI for their survival. Technological progress will become unintelligible to humans, and the TAI would have replaced human-made machines by new forms of physical capital such as robo-factories and specialized compute, which could not be operated without the TAI cognitive input.

Now imagine that, at some finite point $t_{stop}$, such TAI stops working. This could happen, for example, when the TAI undergoes a self-modification allowing it to "wirehead" and achieve unbounded rewards despite not following the goal of utility maximization. Alternatively, the TAI could expect negative utility in the future and decide to optimally switch off at $t_{stop}$ instead.

After the TAI switches off, humanity has to rely on own physical and cognitive labor to re-build capital and recover as much technology as possible, potentially reverting to some primitive technological state amidst rampant violent conflict. Depending on the state of the world at $t_{stop}$, this may or may not be an existential threat.[19]

## 4.3 Social Welfare Function for the TAI

Unfortunately for the prospects of successful TAI alignment, the true functional form of $u(\cdot)$ in (6)–(7) is not known; it is equally difficult to assess whether $u(\cdot)$ should only take total human consumption as input, or individual consumption levels should be aggregated differently. For example, consider three classical cases of social welfare functions that the TAI could potentially maximize:

1. *Benthamite case*: The TAI strives to maximize *aggregate* welfare and, thus, cares for individual utility *and* the number of people. The utility maximization problem of the TAI would then be

$$\max_{\{C(t)\}_{t=0}^{\infty}} N_0 \int_0^{\infty} e^{(n-\rho)t} u(C(t)) M(t) dt, \qquad s.t. \qquad \dot{K} = Y - C - \delta K, \qquad (14)$$

where $n$ is the population growth rate. This is in line with the utility maximization problem of a social planner as it is usually assumed within the framework of Ramsey (1928), Cass (1965), and Koopmans (1965).

2. *Millian case (time-indexed average utilitarianism; Grill, 2023)*: The TAI does not care for the number of people but strives to maximize average utility across individuals. In this case the utility function would be

$$\max_{\{C(t)\}_{t=0}^{\infty}} \int_0^{\infty} e^{-\rho t} u(C(t)) M(t) dt, \qquad s.t. \qquad \dot{K} = Y - C - \delta K. \qquad (15)$$

---

[18]This is akin of zoo animals like in Tegmark's "zookeeper" scenario (Tegmark, 2017).

[19]For example, *immediate* wireheading could cause the TAI to stop working at $t = 0$, before it caused any harm. But that would fall under the "no AI takeover" scenario.

3. *Rawlsian case*: The TAI has an extreme preference for equality and only cares for the least well-off individual. In this case, the utility maximization problem would be

$$\max_{\{C(t)\}_{t=0}^{\infty}} \int_0^{\infty} e^{-\rho t} u(\min\{c_1(t), c_2(t), \ldots c_N(t)\}) M(t) dt, \qquad s.t. \qquad \dot{K} = Y - C - \delta K,$$

(16)

where $N$ is the population size, $i \in [1, 2, \ldots, N]$, and $\sum_{i=1}^{N} c_i(t) = C(t)$. The TAI would have an incentive to redistribute consumption from the rich to the poor until everyone is equally well off.

The differences among these three cases highlight the inter- and intratemporal trade-offs in the maximization of social welfare. All three cases would lead to vastly different outcomes—particularly under sufficiently high TAI optimization power, implying its ability to freely redistribute consumption and control population size. The intertemporal trade-off is clearly visible when comparing the Benthamite and Millian cases. In the Benthamite case, one obtains the "repugnant conclusion" (Parfit, 1986), entailing a maximally large population of people living lives only "barely worth living". By contrast, in the case of the Millian utility function, the TAI would have a preference for keeping population numbers low, which could mean anything from birth control policies to outright killing.

In turn, the intratemporal trade-off is visible when comparing the Millian and the Rawlsian case. Keeping population size constant, both cases can be viewed as polar opposites: in the Millian case, the TAI would disregard the composition of consumption, caring only about its total volume; in the Rawlsian case, the TAI would redistribute until everybody ends up with the same amount of consumption.

All these three cases can be viewed as polar cases of a more general specification of a social welfare function governed by two parameters: the inequality aversion parameter $\theta \in (-\infty, 1]$ and the population size elasticity parameter $\nu \in [0, 1]$, as in

$$\max_{\{C(t)\}_{t=0}^{\infty}} \int_0^{\infty} N(t)^{\nu} e^{-\rho t} u \left( \left[ \sum_{i=1}^{N} c_i(t)^{\theta} \right]^{\frac{1}{\theta}} \right) M(t) dt,$$

$$s.t. \qquad \dot{K} = Y - C - \delta K.$$

(17)

The Benthamite case is obtained by setting $\theta = 1$ and $\nu = 1$; the Millian case by setting $\theta = 1$ and $\nu = 0$; and the Rawlsian case emerges as the limiting case for $\theta \to -\infty$ and $\nu = 0$.

The bottom line of the discussion here is that even if the TAI has a welfare function focusing on human well-being, the parameters of the welfare function will certainly be different from the parameters of *some* individuals and even wide differences between the parameters of the TAI and the *average* parameters across the whole population cannot be ruled out by any means. Thus, even in the optimistic scenario in which the goals of humans and of the TAI are aligned *in principle*, dystopian outcomes for *many* people are still plausible.

18

# 5 How Much Existential Risk Would a Benevolent Social Planner Tolerate?

The main promise of TAI is to accelerate technological progress and economic growth massively, but this hopeful prospect comes at the cost of the risk of human extinction after an AI takeover. To ensure survival, humankind should then either give up on TAI (a scenario that will, in the language of Section 3, materialize with probability $1 - p_1$) or hope that TAI does not take over (with probability $p_1(1 - p_2)$). If any of these scenarios materializes, we would end up in a world with no AI existential risk, but also—due to human decisions bottlenecking the economy—without the massive growth acceleration the technology promises. With that in mind, the most important question in quantifying socially tolerable $p(doom)$ is to weigh the gains of massively accelerated growth after AI takeover against the cost of increased existential risk.

In a first-best world, the decision whether to develop TAI despite the associated existential risk should be informed by an objective comparison of aggregate social welfare obtained in the two scenarios, with and without AI takeover. The benevolent social planner would then optimally steer progress in the direction that yields greater welfare.

In the scenarios in which human extinction is immediate and certain, the trade-off is obvious—humankind is always better off alive. The result is also obvious for the case of zero discounting (cf. MacAskill, 2022): the value of a future in which humans never go extinct is infinite, whereas if extinction happens at any moment in time, that value is finite and hence inferior. On the other hand, with positive discounting, immediate gains from growth acceleration may outweigh the costs of extinction at some point in time $T$, even in scenarios in which eventual extinction is certain. So which path should humanity take?

In what follows, we quantify social welfare under different scenarios and assess the representative individual's willingness to pay for AI safety and alignment research. We assume a standard iso-elastic (constant relative risk aversion, CRRA) utility function of the form

$$u(C(t)) = \frac{C(t)^{1-\theta} - 1}{1 - \theta}, \tag{18}$$

where $C(t)$ is individual consumption at time $t$ and $\theta$ is the coefficient of relative risk aversion, such that the elasticity of intertemporal substitution is $1/\theta$. For this specification, utility is always positive as long as $C(t) > 1$; the logarithmic utility function $u(C(t)) = \log(C(t))$ is obtained as the special case of $\theta \to 1$. We assume $C(t) \geq 1$ for all $t \geq 0$, implying non-negative flow utility and a non-negative extinction hazard rate $m(t)$. After an appropriate normalization, this restriction can be understood as consumption exceeding a minimum subsistence threshold.[20] Ensuring non-negative flow utility is essential in our exercise because we assume that upon death, utility drops to zero.

In what follows, we concentrate on the parameter range $\theta \geq 1$, which includes the most empirically relevant degrees of risk aversion among the human population (see, e.g., Hall, 1988;

---

[20]That is, $C = C_{total}/C_{min} \geq 1$, such that if $C_{total}$ falls below $C_{min}$, the individual dies. Normalization is also helpful in ensuring that the utility function is invariant to the choice of measurement units.

Chetty, 2006; Guvenen, 2006). Specifically, we compare the logarithmic case $\theta = 1$ with the more risk-averse case $\theta > 1$. The key difference between both cases is that for $\theta = 1$ utility is unbounded, whereas for $\theta > 1$ it is bounded above by $1/(\theta - 1) > 0$. This means that growth accelerations are much less valuable with $\theta > 1$, because then even exorbitantly high consumption levels can only generate finite streams of utility (Jones, 2024; Aschenbrenner and Trammell, 2024). Thus, utility gains achieved in "cornucopia" are in expectation much more likely to outweigh the loss of future utility after human extinction if $\theta = 1$ than if $\theta > 1$.

However, it remains unclear how risk aversion aggregates across the human population. Specifically, our individual willingness to accept the risk of death may be different from humanity's willingness to accept the risk of extinction. First, there are evolutionary reasons to believe that individual risk aversion may not represent our attitudes towards aggregate risk correctly: "evolutionarily optimal strategies should be more averse to aggregate risk than to equivalent idiosyncratic risk" (Robson and Samuelson, 2010). Second, specifically in the case of TAI development, the decision whether or not to take the gamble may be taken unilaterally by a selected few individuals, such as the leaders of frontier AI labs, whose past business successes would suggest below-average risk aversion; they could also feel forced to take the gamble due to competitive pressures. Because of these reasons, humanity as a whole may turn out to be much less risk averse than the average or median human, or even risk loving. The case of $\theta = 1$, with unbounded utility, could then be viewed as one in which the social planner internalizes at least part of the social dynamics working in favor of the development of risky TAI, while in the case $\theta > 1$ these dynamics are suppressed.

## 5.1 No AI Takeover Scenario

Without AI takeover, human cognitive work remains the bottleneck of production, and we approximate $X \approx \alpha K$ and $S \approx \gamma AhN$. Therefore $Y \approx F(\alpha K, \gamma AhN)$. In this scenario, long-run growth follows fundamentally from labor-augmenting technical change (Growiec, 2023a,b). Using historical data, we approximate the growth rate $g$ in this scenario as 1.75%.[21] In the absence of existential risk, in this scenario welfare amounts to

$$W_0 = \int_0^\infty e^{-\rho t} \cdot \frac{(c_0 e^{gt})^{1-\theta} - 1}{1 - \theta} \, dt, \tag{19}$$

where $c_0 > 1$ is the exogenous consumption level at $t = 0$.

---

[21]In reality, this growth rate may be lower, because in the absence of TAI, continued growth in the coming decades would require a sustained rate of labor-augmenting technical change. That is, in turn, threatened by unfavorable demographics (such as aging populations), diminishing returns to R&D, and other headwinds to global economic growth such as the need to cut down carbon emissions (Jones, 2002; Gordon, 2016). A secular stagnation is expected, for example, in the OECD long-run GDP forecast until 2060, `https://www.oecd.org/en/data/indicators/real-gdp-long-term-forecast.html`. On the other hand, even narrow AI applications, excluding applications in R&D, can provide a measurable boost to productivity growth, which could counteract these headwinds (Unsal et al., 2024).

## 5.2 AI Takeover Scenarios

With AI takeover at time $t = 0$, human cognitive work can be replaced with the digital software input provided by the TAI. We may approximate $X \approx \alpha K$ and $S \approx A\psi\chi K$. Therefore

$$Y \approx \alpha K F\left(1, \frac{A\psi\chi}{\alpha}\right) \approx \alpha a_K K.$$

In this scenario, there is endogenous growth thanks to the accumulation of compute and the proportional scaling of digital software. The growth rate is an increasing function of the TAI's saving rate (see Prettner, 2019; Lankisch et al., 2019; Growiec, 2023b). Using historical data on the accumulation of compute (Moore's Law), $g^{AI}$ may be about 20%-30% (Davidson, 2021; Growiec, 2022a). However, to ensure that we capture a wide range of plausible outcomes for this case, we use a range of $5 - 40\%$ instead.

### 5.2.1 Well-Aligned and Corrigible TAI

In the case of well-aligned and corrigible TAI, $M(t) = 1$ for all $t$, and therefore total welfare amounts to

$$W_A = \int_0^\infty e^{-\rho t} \cdot \frac{(c_0 e^{g^{AI}t})^{1-\theta} - 1}{1 - \theta} \, dt. \tag{20}$$

As this scenario involves a growth acceleration without incurring any additional costs or risks, it is always strictly preferred to the no-takeover scenario, i.e., we have $W_A > W_0$.

### 5.2.2 One-Off Extinction Risk

For cases with certain human extinction at a future date $T$ (so that $ET = T$), factoring in the existential risk, we obtain

$$W_B = \int_0^T e^{-\rho t} \cdot \frac{(c_0 e^{g^{AI}t})^{1-\theta} - 1}{1 - \theta} \, dt. \tag{21}$$

Now there is a trade-off: the prospect of accelerated growth has to be balanced against the extinction risk (Jones, 2024). Thus, to characterize the trade-off, we compare $W_B$ and $W_0$ and compute the extinction date $T$ for which the benevolent social planner would be indifferent between not developing TAI and developing TAI. We do this using several combinations of the time preference rate $\rho$, the growth rate of the economy in the presence of TAI, $g^{AI}$, and the coefficient of relative risk aversion, $\theta$.

The results are displayed in Table 1. As is intuitive, we observe that the accepted extinction date $T$ is higher with a higher risk aversion ($\theta$), a more patient population (lower $\rho$)—such that the future, and, thus, the time beyond $T$ gets more weight in utility—and a lower growth rate achieved through TAI ($g^{AI}$).

Table 1: Extinction time $T$ (in years from AI takeover), subject to which individuals are indifferent between the scenarios with TAI and without TAI, for varying $\theta$, $g^{AI}$, and $\rho$

| $g^{AI}/\rho$ | $\theta = 1$ | | | | $\theta = 2$ | | | |
|---|---|---|---|---|---|---|---|---|
| | 0.002 | 0.01 | 0.03 | 0.05 | 0.002 | 0.01 | 0.03 | 0.05 |
| 0.05 | 754.47 | 208.68 | 91.8 | 62.63 | 6752.59 | 1238.26 | 403.94 | 243.64 |
| 0.1 | 488.74 | 145.35 | 67.51 | 47.21 | 6623.67 | 1205.72 | 389.07 | 233.12 |
| 0.2 | 331.36 | 103.76 | 50.26 | 35.86 | 6568.34 | 1190.99 | 381.62 | 227.45 |
| 0.3 | 266.57 | 85.45 | 42.23 | 30.45 | 6550.96 | 1186.24 | 379.06 | 225.45 |
| 0.4 | 229.04 | 74.45 | 37.27 | 27.05 | 6542.45 | 1183.89 | 377.82 | 224.41 |

### 5.2.3 Extinction Risk From Misaligned TAI

Next, we consider the probability of TAI misalignment $p_3$. In this case, we compare $W_0$ with $W_D = p_3 \cdot 0 + (1 - p_3)W_A = (1 - p_3)W_A$. This scenario assumes immediate extinction with misaligned TAI ($W_B = 0$). We then ask what the $p(doom)$ is that the benevolent social planner is willing to tolerate to achieve "cornucopia", depending on $p_3$.

Table 2: Values for the probability of immediate misalignment, subject to which individuals are indifferent between the scenarios with TAI and without TAI, for varying $\theta$, $g^{AI}$, and $\rho$

| $g^{AI}/\rho$ | $\theta = 1$ | | | | $\theta = 2$ | | | |
|---|---|---|---|---|---|---|---|---|
| | 0.002 | 0.01 | 0.03 | 0.05 | 0.002 | 0.01 | 0.03 | 0.05 |
| 0.05 | 0.454445 | 0.206246 | 0.0871928 | 0.055282 | 0.00000136 | 0.00000419 | 0.00000546 | 0.00000512 |
| 0.1 | 0.678924 | 0.397439 | 0.195157 | 0.129332 | 0.00000177 | 0.00000580 | 0.00000853 | 0.00000867 |
| 0.2 | 0.823869 | 0.593343 | 0.349124 | 0.247325 | 0.00000197 | 0.00000672 | 0.00001066 | 0.00001151 |
| 0.3 | 0.87865 | 0.693117 | 0.453642 | 0.337154 | 0.00000204 | 0.00000705 | 0.00001150 | 0.00001272 |
| 0.4 | 0.907439 | 0.753577 | 0.529238 | 0.407828 | 0.00000208 | 0.00000722 | 0.00001195 | 0.00001340 |

The results of this calculation are displayed in Table 2. Again, we observe that the risk the social planner would be willing to tolerate is much lower in the case of greater risk aversion. In fact, in the case of $\theta = 2$, almost no extinction risk can be tolerated. By contrast, under moderate risk aversion (log utility, $\theta = 1$), where flow utility is unbounded, the social planner is much more willing to accept such a risk.

In line with intuition, the tolerable extinction risk increases with the growth rate $g^{AI}$ because it allows for much higher consumption in the future, and, thus, a higher lifetime utility under the scenario with TAI. However, in contrast to the previous scenario, more patient individuals would now accept a higher extinction risk. The reason is that extinction would happen immediately with the probability $p_3$, while, if extinction does not occur (with probability $(1 - p_3)$), growth at a high rate of $g^{AI}$ sets in. If individuals have a lower $\rho$, the discounted lifetime utility of the future consumption that can be achieved with the higher growth rate is also higher. As a consequence, more patient individuals would tolerate a higher immediate extinction risk.

### 5.2.4 Extinction Risk From Non-Corrigible TAI

In the next scenario, we consider the probability of non-corrigible TAI, $p_4$. In this case, the benevolent social planner would compare $W_0$ with $W_E = p_3 \cdot 0 + (1 - p_3)p_4 W_B + (1 - p_3)(1 - p_4)W_A$, where we assume immediate extinction with misaligned TAI ($W_B = 0$), as well as extinction at

time $T$ with initially aligned but non-corrigible TAI ($W_B > 0$, depending on $T$). The probability of humanity's long-term survival is now $M_\infty = (1 - p_3)(1 - p_4)$. The question is which $p(doom)$ the social planner would be willing to tolerate, depending on $p_3$, $p_4$, and $T$.

The results now come in three different parts. In Table 3, Panel A, we display the values of the probability of immediate misalignment ($p_3$) subject to which the social planner would be indifferent between the scenarios with TAI and without TAI, holding the probability of non-corrigibility ($p_4$) and the time after which the non-corrigibility of TAI leads to doom ($T$) constant. In Panel B, we display the values of the probability of non-corrigibility ($p_4$) subject to which the social planner would be indifferent between the scenarios with TAI and without TAI, holding the probability of immediate misalignment ($p_3$) and the time after which the non-corrigibility of TAI leads to doom ($T$) constant. Finally, in Panel C, we display the values of the time after which the non-corrigibility of TAI leads to doom ($T$), subject to which the social planner would be indifferent between the scenarios with TAI and without TAI, holding the probability of immediate misalignment ($p_3$) and the probability of non-corrigibility ($p_4$) constant.

We again observe that the social planner is much less willing to tolerate extinction risk if the coefficient of relative risk aversion $\theta$ is higher. As far as the growth rate $g^{AI}$ is concerned, a higher growth rate leads to higher risks that individuals would accept, either in terms of higher probabilities of immediate misalignment ($p_3$) and non-corrigibility ($p_4$), or in terms of a smaller number of years after which non-corrigibility leads to doom ($T$). Again, these results are very much in line with economic intuition.

### 5.2.5 Mounting Extinction Risk

Finally, we consider mounting extinction risk, that is, risk, which builds up gradually in the background. We assume that, other things equal, the more advanced TAI becomes and the higher levels of output it provides, the higher is the extinction risk. We capture this by assuming that the extinction risk rises with consumption such that

$$W_C = \int_0^\infty e^{-\rho t - \int_0^t m(s)ds} \cdot \frac{(c_0 e^{g^{AI}t})^{1-\theta} - 1}{1 - \theta} \, dt, \tag{22}$$

where $m(s) = \log C(s)^\varepsilon = \varepsilon \cdot [\log(c_0) + g^{AI} \cdot s]$, in which $\varepsilon$ measures the strength of the relationship between consumption and existential risk. In this specification, which is by construction more conservative than that of Aschenbrenner and Trammell (2024), the instantaneous extinction hazard rate is proportional to the *log* of consumption, and thus as consumption grows exponentially over time, the hazard rate increases only arithmetically. Eventual extinction is now certain ($M_\infty = 0$), and humanity's expected lifespan after AI takeover is a decreasing function

Table 3: Values for the probability of immediate misalignment (Panel A), the probability of the AI being non-corrigible (Panel B), and the time $T$ at which the AI can become non-corrigible (Panel C) subject to which the social planner is indifferent between the scenarios with TAI and without TAI.

| | Panel A: $p_4 = 0.00003$ and $T = 50$ | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | $\theta = 1$ | | | | $\theta = 2$ | | | |
| $g^{AI}/\rho$ | 0.002 | 0.01 | 0.03 | 0.05 | 0.002 | 0.01 | 0.03 | 0.05 |
| 0.05 | 0.454429 | 0.206229 | 0.0871855 | 0.0552791 | - | - | - | 0.0000026596 |
| 0.1 | 0.678915 | 0.397425 | 0.195149 | 0.129329 | - | - | 0.0000018340 | 0.0000062058 |
| 0.2 | 0.823864 | 0.593334 | 0.349117 | 0.247322 | - | - | 0.0000039688 | 0.0000090426 |
| 0.3 | 0.878647 | 0.693109 | 0.453636 | 0.337151 | - | - | 0.0000048098 | 0.0000102584 |
| 0.4 | 0.907437 | 0.753571 | 0.529232 | 0.407825 | - | - | 0.0000052596 | 0.0000109339 |

| | Panel B: $p_3 = 0.00003$ and $T = 50$ | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | $\theta = 1$ | | | | $\theta = 2$ | | | |
| $g^{AI}/\rho$ | 0.002 | 0.01 | 0.03 | 0.05 | 0.002 | 0.01 | 0.03 | 0.05 |
| 0.05 | 0.469403 | 0.293447 | 0.325211 | 0.5551 | - | - | - | - |
| 0.1 | 0.693265 | 0.528045 | 0.645486 | $> 1$ | - | - | - | - |
| 0.2 | 0.835111 | 0.738226 | 0.994075 | $> 1$ | - | - | - | - |
| 0.3 | 0.888181 | 0.835321 | $> 1$ | $> 1$ | - | - | - | - |
| 0.4 | 0.915953 | 0.89125 | $> 1$ | $> 1$ | - | - | - | - |

| | Panel C: $p_3 = p_4 = 0.3$ | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | $\theta = 1$ | | | | $\theta = 2$ | | | |
| $g^{AI}/\rho$ | 0.002 | 0.01 | 0.03 | 0.05 | 0.002 | 0.01 | 0.03 | 0.05 |
| 0.05 | 355.307 | 9308.72 | 2533.34 | 323.316 | 35750.9 | 4250.81 | 1785.45 | 418.062 |
| 0.1 | - | 123.461 | 1180.94 | 817.496 | 35750.8 | 4250.79 | 1785.44 | 418.059 |
| 0.2 | - | - | 67.5608 | 902.427 | 35750.8 | 4250.79 | 1785.43 | 418.056 |
| 0.3 | - | - | 18.1997 | 46.815 | 35750.8 | 4250.79 | 1785.43 | 418.055 |
| 0.4 | - | - | - | 20.6203 | 35750.8 | 4250.78 | 1785.43 | 418.055 |

Note that "-" means that the implied probability or the implied time $T$ would be negative. Thus, the no-TAI scenario is always preferred. In turn, "$> 1$" means that the implied probability would be above 1. Thus, the TAI scenario is always preferred.

of $g^{AI}$. Specifically, under the simplifying assumption $c_0 = 1$ this lifespan is equal to[22]

$$ET = \sqrt{\frac{\pi}{2\varepsilon g^{AI}}}.$$

In a world run by TAI that does not mitigate existential risks, growth accelerations increase the extinction hazard rate.

In Table 4, we display the values of $\varepsilon$ at which the social planner would be indifferent between the scenarios with TAI and without TAI. We observe that with a higher risk aversion, the planner would never choose the scenario with TAI. If risk aversion is close to $\theta \approx 1$, though,

---

[22]For $c_0 > 1$, humanity's expected lifespan is

$$ET = \sqrt{\frac{\pi}{2\varepsilon g^{AI}}} e^{\frac{\varepsilon \cdot (\log(c_0))^2}{2g^{AI}}} \left(1 - \text{erf}\left(\frac{\sqrt{\varepsilon}\log(c_0)}{\sqrt{2g^{AI}}}\right)\right).$$

faster economic growth in terms of $g^{AI}$ makes the TAI scenario more attractive, while a higher time preference rate $\rho$ is a double-edged sword. On the one hand, it implies that the future risk of extinction is discounted more heavily; on the other hand, also the benefits of TAI through exponential consumption growth are discounted more heavily. For comparatively low $g^{AI}$, this implies that no solution can be found in terms of $\varepsilon$ for which the TAI scenario is preferred.

Table 4: Values for $\varepsilon$ in the case of mounting extinction risk, at which individuals are indifferent between the scenarios with TAI and without TAI, for $\theta = 1.0001$, $\theta = 2$, and varying $g^{AI}$ and $\rho$

| | $\theta = 1.0001$ | | | | $\theta = 2$ | | | |
|---|---|---|---|---|---|---|---|---|
| $g^{AI}/\rho$ | 0.002 | 0.01 | 0.03 | 0.05 | 0.002 | 0.01 | 0.03 | 0.05 |
| 0.05 | 0.000 027 89 | 0.000 124 51 | - | - | - | - | - | - |
| 0.1 | 0.000 043 92 | 0.000 223 22 | 0.000 419 39 | - | - | - | - | - |
| 0.2 | 0.000 057 48 | 0.000 323 71 | 0.000 691 51 | 0.000 888 42 | - | - | - | - |
| 0.3 | 0.000 064 22 | 0.000 380 09 | 0.000 865 70 | 0.001 155 22 | - | - | - | - |
| 0.4 | 0.000 068 47 | 0.000 418 17 | 0.000 991 82 | 0.001 356 91 | - | - | - | - |

Note that "-" means that no value could be found numerically. In that case, no TAI is always preferred. We used $\theta = 1.0001$ as a representation for the logarithmic case, because the algorithm did not converge otherwise.

## 5.3 Willingness to Pay for Doom Avoidance

We now quantify the benevolent social planner's willingness to pay ($WTP$) to avoid the existential risk stemming from TAI (cf. Shulman and Thornley, 2024; Jones, 2025). We do this by computing the equivalent variation ($EV$), that is, the factor by which the consumption level of an individual living in cornucopia (i.e., in a world with well-aligned and corrigible TAI) would have to be reduced to make them equally well off as an individual living in a world in which the TAI could be potentially misaligned or non-corrigible. The fraction of consumption in each period that the representative individual would be willing to forgo to guarantee an indefinite prosperous life in cornucopia instead of facing extinction risk is then given by $1 - EV$. $EV$ is defined as

$$EV = e^{-(W_0 - W_B) \cdot \rho}, \tag{23}$$

such that the willingness to pay to avoid doom in the various scenarios amounts to

$$WTP = (1 - EV) \cdot c_t \tag{24}$$

in each period $t$. In the following scenarios, we only consider the case of $\theta = 1$ because for $\theta = 2$, the social planner would either prefer not to build TAI at all, or would be willing to forgo almost all consumption to avoid existential risk. Restricting attention to the case of $\theta = 1$ again ensures that we are on the conservative side of assessing $EV$ and the associated $WTP$—i.e., we provide a lower bound for the actual $WTP$—because $\theta$ tends to be higher than unity according to empirical estimates.

Table 5, Panel A contains the $EV$ for avoiding AI-driven extinction after 100 years. The results clearly show a very high $WTP$: the social planner would be willing to forgo a large share of consumption to avoid extinction in 100 years, particularly if it is patient. If it is impatient, however, the prospect of extinction in 100 years looses its grip and the $EV$ rises to values above

Table 5: Equivalent variation (EV) for $\theta = 1$, varying $g^{AI}$ and $\rho$, and different scenarios for $T$, $p_3$, $p_4$, and $\varepsilon$

| | Panel A: EV for $\theta = 1$, $T = 100$, and varying $g^{AI}$ and $\rho$ | | | |
|---|---|---|---|---|
| $g^{AI}/\rho$ | 0.002 | 0.01 | 0.03 | 0.05 |
| 0.05 | $3.22 \times 10^{-15}$ | 0.000 482 56 | 0.419 993 | 0.893 228 |
| 0.1 | $6.927\,79 \times 10^{-26}$ | $1.218\,63 \times 10^{-5}$ | 0.301 365 | 0.857 837 |
| 0.2 | $3.209\,08 \times 10^{-47}$ | $7.771\,61 \times 10^{-9}$ | 0.155 166 | 0.791 206 |
| 0.3 | $1.4865 \times 10^{-68}$ | $4.956\,22 \times 10^{-12}$ | 0.079 891 4 | 0.729 751 |
| 0.4 | $6.885\,73 \times 10^{-90}$ | $3.160\,75 \times 10^{-15}$ | 0.041 134 2 | 0.673 069 |

| | Panel B: EV for $\theta = 1$, $p_3 = 0.1$, and varying $g^{AI}$ and $\rho$ | | | |
|---|---|---|---|---|
| $g^{AI}/\rho$ | 0.002 | 0.01 | 0.03 | 0.05 |
| 0.05 | 0.027 993 3 | 0.206 844 | 0.288 674 | 0.308 575 |
| 0.1 | 0.002 297 83 | 0.125 457 | 0.244 357 | 0.279 21 |
| 0.2 | $1.548\,27 \times 10^{-5}$ | 0.046 153 1 | 0.175 09 | 0.228 598 |
| 0.3 | $1.043\,21 \times 10^{-7}$ | 0.016 978 8 | 0.125 457 | 0.187 16 |
| 0.4 | $7.029\,11 \times 10^{-10}$ | 0.006 246 15 | 0.089 894 | 0.153 234 |

| | Panel C: EV for $\theta = 1$, $p_3 = p_4 = 0.1$, $T = 50$, and varying $g^{AI}$ and $\rho$ | | | |
|---|---|---|---|---|
| $g^{AI}/\rho$ | 0.002 | 0.01 | 0.03 | 0.05 |
| 0.05 | $1.24 \times 10^{-3}$ | 0.076 348 4 | 0.213 917 | 0.277 725 |
| 0.1 | $1.085\,58 \times 10^{-5}$ | 0.030 750 3 | 0.166 542 | 0.244 882 |
| 0.2 | $8.298\,68 \times 10^{-10}$ | 0.004 988 24 | 0.100 944 | 0.190 388 |
| 0.3 | $6.343\,89 \times 10^{-14}$ | 0.000 809 18 | 0.061 184 | 0.148 02 |
| 0.4 | $4.849\,55 \times 10^{-18}$ | 0.000 131 263 | 0.037 084 7 | 0.115 081 |

| | Panel D: EV for $\theta = 1$, $\varepsilon$ as in Table 4, and varying $g^{AI}$ and $\rho$ | | | |
|---|---|---|---|---|
| $g^{AI}/\rho$ | 0.002 | 0.01 | 0.03 | 0.05 |
| 0.05 | $9.41 \times 10^{-8}$ | 0.038 994 2 | - | - |
| 0.1 | $3.431\,94 \times 10^{-16}$ | 0.000 266 125 | 0.064 185 9 | - |
| 0.2 | $6.906\,44 \times 10^{-40}$ | $1.258\,46 \times 10^{-8}$ | 0.002 305 66 | 0.002 557 01 |
| 0.3 | $4.783\,92 \times 10^{-61}$ | $6.0691 \times 10^{-13}$ | $8.3005 \times 10^{-5}$ | 0.003 551 61 |
| 0.4 | $5.355\,08 \times 10^{-82}$ | $2.985\,82 \times 10^{-17}$ | $2.994\,86 \times 10^{-6}$ | 0.000 483 049 |

65% (for $\rho = 0.05$). Panel B contains the $EV$ for avoiding the risk of immediate extinction if it amounts to $p_3 = 0.1$. Now the results are, of course, much less sensitive to the time preference rate because the risk is immediate. Again, the results show a high willingness to pay (70% of yearly consumption and more) to reduce the risk from 10% to zero. Panel C contains the $EV$ for avoiding the risk of immediate extinction if it amounts to $p_3 = 0.1$ *and* the risk of non-corrigibility if it amounts to $p_4 = 0.1$ and leads to extinction after $T = 50$ years. As compared with the results in Panel B, $EV$ is typically smaller, which is not surprising because of the additional risk of non-corrigibility that individuals can pay to avoid here. In addition, $EV$ increases with the time preference rate $\rho$ towards the results from Panel B, which is again intuitive because non-corrigibility leads to extinction after 50 years and the higher the time preference rate, the

lower is the effect of the time after $T = 50$ on lifetime utility. Finally, in Panel D, we compute the $EV$ associated with mounting extinction risk. We do this not for a pre-specified value of $\varepsilon$, because in these cases the numerical method only works well for combinations of $g^{AI}$ and $\rho$ that lead to values of $\varepsilon$ that are rather close to the value we assumed that the social planner wants to eliminate. By contrast, we always assume that the value of $\varepsilon$ that the social planner wants to avoid is the same as the corresponding entry in Table 4. We observe that the willingness to avoid the associated risks (which are already close to zero to start with) is high and amounts to almost all of consumption. Again, $EV$ increases with the time preference rate $\rho$ because more impatience means that the social planner puts less weight on the future when the extinction risk has become very high already. As far as the growth rate $g^{AI}$ is concerned, $EV$ decreases and the willingness to pay increases with that rate because the stakes become higher if growth is faster.

## 5.4 Discussion

Let us now put these numbers in perspective. How much existential risk from TAI would the benevolent social planner tolerate, and how does that compare to the $p(doom)$ figures mentioned in the popular discourse?

In the baseline scenario, featuring a realistic degree of risk aversion ($\theta = 2$), a moderate 3% annual discount rate, and the optimistic assumption that TAI would accelerate economic growth to 30% per annum, we find that TAI should be developed:

- in the case where human extinction is certain—only if extinction happens no earlier than in 379 years from AI takeover,

- in the case of a one-off extinction risk occurring immediately at AI takeover—only if $p(doom)$ is below 0.00001 (one in a hundred thousand),

- in the case of a 30% risk of extinction immediately upon AI takeover and a 30% conditional risk of extinction later—only if that later extinction hazard materializes no earlier than after 1785 years,

- in the case TAI will bring continuous extinction risk, gradually increasing log-linearly with humanity's per capita consumption—never.

Given that it would be close to impossible to halt AI development, how much should society pay to avoid AI doom? To estimate this amount, we have compared two scenarios: one with a risky TAI, and one in which "cornucopia" is certain, but achieved at the cost of a fraction of consumption spent each year, from $t = 0$ until infinity, on existential risk mitigation.

We find that with $\theta = 2$, the benevolent social planner would be willing to pay almost all of consumption to avoid human extinction. But we find very high numbers even for the case $\theta = 1$, in which flow utility is allowed to be arbitrarily high:

- the acceptable price for avoiding certain human extinction in 100 years from AI takeover is 92% of total consumption each year ($EV = 0.080$),[23]

- the acceptable price for avoiding a 10% chance of human extinction immediately upon AI takeover is 87.5% of total consumption each year ($EV = 0.125$),

- the acceptable price for avoiding a 10% chance of human extinction immediately upon AI takeover and a 10% conditional chance of human extinction 50 years later is 93.9% of total consumption each year ($EV = 0.061$),

- in the case in which TAI brings continuous extinction risk that increases log-linearly with humanity's per capita consumption—the acceptable price is almost all consumption each year ($EV = 8 \times 10^{-5}$).

Recall that these are not one-off costs to avoid AI doom, but costs that must be borne each year, repeatedly, throughout the entire future.

AI experts and technology leaders have repeatedly mentioned $p(doom)$ guesstimates for the next 25-75 years that exceed 10%. Despite that, spending on AI safety and AI alignment research is close to zero percent of global consumption or GDP; it is tiny even when compared to research spending of the AI industry alone. According to Hilton (2022), in 2022 there were only about 400 people worldwide working on AI alignment; "around \$50 million was spent on reducing catastrophic risks from AI in 2020—while billions were spent advancing AI capabilities". According to Bengio et al. (2024), "only an estimated 1–3% of AI publications are on safety". This means that underinvestment in AI safety and AI alignment research is colossal. Thus, humanity is exposed to excessive amounts of existential risk.

## 5.5 Extensions

Our numerical exercises are also informative for certain additional cases.

First, instead of the Millian social welfare function, one could consider the Benthamite case. However, under exogenous population growth, this boils down to simply changing the social discount rate from $\rho$ to $\rho - n$. Even more generally, for any intermediate case between the Millian and the Benthamite social welfare function, one could replace the discount rate with $\rho - \nu n$, where $\nu \in (0, 1)$. Then, as long as the population growth rate is not affected by our choice whether or not to adopt risky TAI, our numerical results go through, pending only the adjustment of the discount rate.

Second, one could add background extinction risk, understood as the risk of extinction from other causes than TAI, such as a deadly pandemic, a nuclear war followed by nuclear winter, an impact of a large asteroid, a gamma ray burst in the vicinity of the solar system, or a supervolcano eruption (Ord, 2020). However, again, under a constant extinction hazard rate

---

[23]Using the latest data of the Federal Reserve Bank of St. Louis on US consumption expenditure (which are for the year 2023), and assuming that mitigation costs would be shared globally in proportion to countries' aggregate consumption, specifically for the case of the USA this amounts to 64% of GDP or 17.3 trillion US dollar.

$m > 0$, this boils down to replacing the social discount rate $\rho$ with $\rho + m$. Then, as long as our choice whether or not to adopt risky TAI does not affect $m$, our results still go through.

Of course, one could argue that TAI may reduce the background extinction hazard rate $m$. For example, it could design vaccines and antidotes, enforce nuclear disarmament, or allow humanity to become multiplanetary. This would be, in fact, equivalent to the Jones (2024) exercise where "AI might create cures for diseases and reduce mortality more generally." Such a change would, under any parametrization, favor the gamble with risky TAI in comparison with the no-TAI baseline. However, it is not clear *a priori* whether TAI would decrease or increase the background extinction hazard rate $m$. On the flip side, it could also increase $m$ by, for example, designing deadly pathogens or facilitating nuclear escalation.

More importantly, the extent of existential risk posed by misaligned or non-corrigible TAI appears to be several orders of magnitude greater than the background hazard rate $m$ that could potentially be reduced by TAI. Factoring in only the natural extinction risks, humanity's expected lifespan ($ET = \int_0^\infty e^{-mt} dt = 1/m$) is probably at least a few hundred thousand years, or even millions of years (Ord, 2020). Considering also the risks due to nuclear or biological weapons reduces $ET$ by perhaps one or two orders of magnitude, depending on one's assessment of the likelihood that any given nuclear war or engineered pandemic could escalate to extinction-level proportions. By contrast, introduction of TAI—without the warranty of its alignment and corrigibility—is likely to reduce humanity's expected lifespan to the range of mere decades.

## 6 Conclusions

We have characterized and analyzed an extensive set of scenarios that may arise with the development of transformative artificial intelligence (TAI). In this context, TAI presents itself as a double-edged sword. On one hand, it holds the promise of achieving a technological singularity, which could lead to unprecedented advancements in human capabilities and societal well-being. On the other hand, it also serves as an important source of existential risk, posing potentially catastrophic threats to humanity. From our theoretical analysis, we derive several important lessons: (i) many plausible scenarios ultimately lead to "AI doom," whereas fewer scenarios point toward a "post-scarcity" future; (ii) AI doom may not necessarily coincide with the moment of AI takeover but could occur later; (iii) even in cases in which TAI prioritizes human well-being in general, large parts of the population may face negative consequences; and (iv) investing in AI safety and alignment research is crucial to mitigating these risks. Our numerical analysis provides a framework to quantify an acceptable level of existential risk, balancing its consequences against the potential acceleration of growth in human consumption and overall well-being. Moreover, our results allow us to estimate how much a social planner would be willing to pay to avoid catastrophic AI outcomes.

Our findings underscore that the results are extremely sensitive to humanity's level of risk aversion. With a realistic degree of risk aversion, a social planner would likely avoid the development of TAI in nearly all scenarios unless TAI is deemed "almost" completely safe. But even under more conservative parameter specifications, such as the logarithmic case in which

flow utility can be unbounded, it is clear that society is currently significantly underinvesting in efforts to reduce existential risks, i.e., in AI safety and AI alignment research. Our results show that a social planner would be willing to forgo a much larger portion of current consumption to ensure that these risks are addressed properly. The key policy implication of our research is that immediate and substantial increases in investment are needed in AI safety and AI alignment research. To ensure that we can minimize existential risks and maximize the potential benefits of TAI, these challenges must effectively be resolved before AI takeover occurs.

# References

Acemoglu, D. (2024). The Simple Macroeconomics of AI. Working paper 32487, National Bureau of Economic Research.

Acemoglu, D. and Restrepo, P. (2018). The Race Between Man and Machine: Implications of Technology for Growth, Factor Shares and Employment. *American Economic Review*, 108:1488–1542.

Agüera y Arcas, B. and Norvig, P. (2023). Artificial General Intelligence Is Already Here. *Noema*, pages October 10, 2023.

Allyn-Feuer, A. and Sanders, T. (2023). Transformative AGI by 2043 is <1% likely. arxiv.org: `https://doi.org/10.48550/arXiv.2306.02519`.

Amodei, D. (2024). Machines of Loving Grace: How AI Could Transform the World for the Better. Essay, Anthropic.

Aschenbrenner, L. (2020). Existential Risk and Growth. Global Priorities Institute WP 6-2020, University of Oxford.

Aschenbrenner, L. (2024). Situational Awareness: The Decade Ahead. Technical report, For Our Posterity.

Aschenbrenner, L. and Trammell, P. (2024). Existential Risk and Growth. Technical report, University of Oxford.

Bengio, Y., Hinton, G., Yao, A., et al. (2024). Managing Extreme AI Risks Amid Rapid Progress. *Science*, 384:842–845.

Bostrom, N. (2014). *Superintelligence: Paths, Dangers, Strategies*. Oxford University Press.

Cass, D. (1965). Optimum Growth in an Aggregative Model of Capital Accumulation. *The Review of Economic Studies*, Vol. 32(No. 3):233–240.

Chetty, R. (2006). A New Method of Estimating Risk Aversion. *American Economic Review*, 96:1821–1834.

Chollet, F. (2019). On the Measure of Intelligence. Working paper, Google, Inc.

Cotra, A. (2022a). Two-Year Update on My Personal AI Timelines. Technical report, AI alignment forum.

Cotra, A. (2022b). Without Specific Countermeasures, The Easiest Path to Transformative AI Likely Leads to AI Takeover. AI Alignment Forum, July 18, 2022, Open Philanthropy.

Davidson, T. (2021). Could Advanced AI Drive Explosive Economic Growth? Technical report, Open Philanthropy.

Davidson, T. (2023). What a Compute-Centric Framework Says About Takeoff Speeds. Technical report, Open Philanthropy.

Dung, L. (2024). The Argument for Near-Term Human Disempowerment Through AI. *AI and Society*, page forthcoming.

Epoch AI (2023). Key Trends and Figures in Machine Learning. Accessed: 2024-07-22.

Epoch AI (2024). Direct Approach Interactive Model. Technical report, Epoch AI.

Field, S. (2025). Why Do Experts Disagree on Existential Risk and $p(doom)$? A Survey of AI Experts. arxiv:2502.14870, Cambridge University, UK.

Gordon, R. J. (2016). *The Rise and Fall of American Growth: The U.S. Standard of Living since the Civil War*. Princeton University Press.

Grace, K., Stewart, H., Sandkühler, J., Thomas, S., Weinstein-Raun, B., and Brauner, J. (2024). Thousands of AI Authors on the Future of AI. Technical report, ArXiv:2401.02843.

Greenblatt, R., Denison, C., Wright, B., Roger, F., MacDiarmid, M., Marks, S., Treutlein, J., Belonax, T., Chen, J., Duvenaud, D., Khan, A., Michael, J., Mindermann, S., Perez, E., Petrini, L., Uesato, J., Kaplan, J., Shlegeris, B., Bowman, S. R., and Hubinger, E. (2024). Alignment Faking in Large Language Models. Working Paper No. 2412.14093, arXiv.

Grill, K. (2023). The Sum of Averages: An Egyptology-Proof Average View. *Utilitas*, 35:103–118.

Growiec, J. (2022a). *Accelerating Economic Growth: Lessons from 200 000 Years of Technological Progress and Human Development*. Springer.

Growiec, J. (2022b). Automation, Partial and Full. *Macroeconomic Dynamics*, 26:1731–1755.

Growiec, J. (2023a). Industry 4.0? Framing the Digital Revolution and Its Long-Run Growth Consequences. *Gospodarka Narodowa. The Polish Journal of Economics*, 316:1–16.

Growiec, J. (2023b). What Will Drive Global Economic Growth in the Digital Age? *Studies in Nonlinear Dynamics and Econometrics*, 27:335–354.

Growiec, J., Jabłońska, J., and Parteka, A. (2024). Hardware and Software: A New Perspective on the Past and Future of Economic Growth. Working Paper, Brookings Institution.

Guvenen, F. (2006). Reconciling Conflicting Evidence on the Elasticity of Intertemporal Substitution: A Macroeconomic Perspective. *Journal of Monetary Economics*, Vol. 53:1451–1472.

Hadfield-Menell, D. (2021). The Principal-Agent Alignment Problem in Artificial Intelligence. Technical report ucb/eecs-2021-207, University of California at Berkeley.

Hall, R. (1988). Intertemporal Substitution in Consumption. *Journal of Political Economy*, 96:339–357.

Hémous, D. and Olsen, M. (2022). The Rise of the Machines: Automation, Horizontal Innovation, and Income Inequality. *American Economic Journal: Macroeconomics*, 14(1):179–223.

Hilton, B. (2022). Preventing an AI-Related Catastrophe. Updated in July 2024, 80 000 Hours.

Jones, C. I. (2002). Sources of U.S. Economic Growth in a World of Ideas. *American Economic Review*, 92:220–239.

Jones, C. I. (2024). The AI Dilemma: Growth versus Existential Risk. *AER: Insights*, 6:575–590.

Jones, C. I. (2025). How Much Should We Spend to Reduce AI's Existential Risk. Unpublished, Stanford GSB.

Karnofsky, H. (2016). Some Background on Our Views Regarding Advanced Artificial Intelligence. Technical report, Open Philanthropy.

Koopmans, T. C. (1965). On the Concept of Optimal Economic Growth. In *The Econometric Approach to Development Planning*. Amsterdam: North Holland.

Korinek, A. and Stiglitz, J. (2019). Artificial Intelligence and Its Implications for Income Distribution and Unemployment. In Agrawal, A., Gans, J. S., and Goldfarb, A., editors, *The Economics of Artificial Intelligence: An Agenda*. University of Chicago Press.

Kurzweil, R. (2005). *The Singularity is Near*. New York: Penguin.

Lankisch, C., Prettner, K., and Prskawetz, A. (2019). How Can Robots Affect Wage Inequality? *Economic Modelling*, 81(September 2019):161–169.

MacAskill, W. (2022). *What We Owe the Future: A Million-Year View*. Basic Books.

Manheim, D. and Garrabrant, S. (2019). Categorizing Variants of Goodhart's Law. arxiv:1803.04585, Machine Intelligence Research Institute.

Nordhaus, W. D. (2021). Are We Approaching an Economic Singularity? Information Technology and the Future of Economic Growth. *American Economic Journal: Macroeconomics*, 13(1):299–332.

Ord, T. (2020). *The Precipice: Existential Risk and the Future of Humanity*. Hachette.

Parfit, D. (1986). *Reasons and Persons*. Oxford University Press.

Peretto, P. and Seater, J. (2013). Factor-Eliminating Technical Change. *Journal of Monetary Economics*, 60:459–473.

Pratt, G. A. (2015). Is a Cambrian Explosion Coming for Robotics? *Journal of Economic Perspectives*, Vol. 29(No. 3):51–60.

Prettner, K. (2019). A Note on the Implications of Automation for Economic Growth and the Labor Share. *Macroeconomic Dynamics*, 23(3):1294–1301.

Prettner, K. and Strulik, H. (2020). Innovation, Automation, and Inequality: Policy Challenges in the Race Against the Machine. *Journal of Monetary Economics*, 116:249–265.

Ramsey, F. P. (1928). A Mathematical Theory of Saving. *Economic Journal*, 38:543–559.

Robson, A. J. and Samuelson, L. (2010). The Evolutionary Foundations of Preferences. In Benhabib, J., Bisin, A., and Jackson, M., editors, *The Social Economics Handbook*. Elsevier.

Roodman, D. (2020). On the probability distribution of long-term changes in the growth rate of the global economy: An outside view. Technical report, Open Philanthropy.

Russell, S. J. (2019). *Human Compatible: Artificial Intelligence and the Problem of Control.* Viking.

Shulman, C. and Thornley, E. (2024). How Much Should Governments Pay to Prevent Catastrophes? Longtermism's Limited Role. In Barrett, J., Greaves, H., and Thorstad, D., editors, *Essays on Longtermism*. Oxford University Press.

Soares, N., Fallenstein, B., Yudkowsky, E., and Armstrong, S. (2015). Corrigibility. Paper presented at AAAI Workshops: Workshops at the Twenty-Ninth AAAI Conference on Artificial Intelligence, Austin, TX, January 25–26, 2015, Machine Intelligence Research Institute.

Steigum, E. (2011). Robotics and Growth. In de la Grandville, O., editor, *Frontiers of Economics and Globalization: Economic Growth and Development*, pages 543–557. Emerald Group. Bingley, UK.

Tegmark, M. (2017). *Life 3.0: Being Human in the Age of Artificial Intelligence*. New York: Knopf.

Tegmark, M. and Omohundro, S. (2023). Provably Safe Systems: The Only Path to Controllable AGI. arxiv:2309.01933, MIT.

Trammell, P. (2021). Existential Risk and Exogenous Growth. Technical report, Global Priorities Institute, University of Oxford.

Trammell, P. and Korinek, A. (2020). Economic Growth Under Transformative AI. Global Priorities Institute WP 8-2020, University of Oxford.

Unsal, F., Koźluk, T., Filippucci, F., Gal, P., and Schief, M. (2024). Miracle or Myth? Assessing the Macroeconomic Productivity Gains from Artificial Intelligence. OECD AI Working Paper No. 29, OECD.

Yampolskiy, R. V. (2024). *AI: Unexplainable, Unpredictable, Uncontrollable.* Chapman & Hall / CRC Artificial Intelligence and Robotics Series.

Yudkowsky, E. (2013). Intelligence Explosion Microeconomics. Technical report 2013-1, Machine Intelligence Research Institute.

Yudkowsky, E. (2022). MIRI announces new "Death With Dignity" strategy. Technical report, Less Wrong, 2 April 2022.

Zhuang, S. and Hadfield-Menell, D. (2021). Consequences of Misaligned AI. arxiv:2302.02083, University of California at Berkeley.