

# Regulating Transformative Technologies

Daron Acemoglu      Todd Lensman\*

January 8, 2024

## Abstract

Transformative technologies like generative AI promise to accelerate productivity growth across many sectors, but they also present new risks from potential misuse. We develop a multi-sector technology adoption model to study the optimal regulation of transformative technologies when society can learn about these risks over time. Socially optimal adoption is gradual and typically convex. If social damages are large and proportional to the new technology's productivity, a higher growth rate paradoxically leads to slower optimal adoption. Equilibrium adoption is inefficient when firms do not internalize all social damages, and sector-independent regulation is helpful but generally not sufficient to restore optimality.

**JEL Classification:** H21, O33, O41

**Keywords:** AI, disasters, economic growth, regulation, technology adoption

Recent breakneck advances in (generative) artificial intelligence have simultaneously raised hopes of productivity gains in many sectors and fears that this technology will be used for nefarious purposes, even posing an existential risk comparable to nuclear war.<sup>1</sup> Some experts have called to slow down or pause the development and adoption of AI technologies,<sup>2</sup> partly because a slower rollout might provide time to identify danger areas and craft appropriate regulations. However, there is little economic analysis of these issues, and it is unclear whether slowing the development and adoption of a promising, transformative technology ever makes sense.

---

\*Acemoglu: Massachusetts Institute of Technology, Department of Economics (email: [daron@mit.edu](mailto:daron@mit.edu)); Lensman: Massachusetts Institute of Technology, Department of Economics (email: [tlensman@mit.edu](mailto:tlensman@mit.edu)). We thank Glen Weyl for several useful discussions; Joshua Gans, Chad Jones, the editor Peter Klenow, and three anonymous referees for comments; and the Hewlett Foundation and the National Science Foundation for financial support.

<sup>1</sup><https://www.nytimes.com/2023/05/30/technology/ai-threat-warning.html>

<sup>2</sup><https://futureoflife.org/open-letter/pause-giant-ai-experiments/>

In this paper, we develop a framework to provide a first set of insights on these questions. We consider a multi-sector economy that initially uses an old technology but can switch to a new, transformative technology. This technology is *transformative* both because it enables a higher growth rate of output, and because it is general-purpose and can be adopted across all sectors of the economy. It also poses new risks. We model these by assuming that there is a positive probability of a *disaster*, meaning that the technology will turn out to have many harmful uses. If a disaster is realized, some of the sectors that had started using the new technology may not be able to switch away from it, despite the social damages. Whether there will be a disaster is initially unknown, and society can learn about it over time. Critically, we assume that the greater are the new technology’s capabilities, the more damaging it will be when used for harmful purposes.<sup>3</sup>

In this environment, we study (socially) optimal and equilibrium adoption decisions. We first show that it is optimal to adopt the new technology gradually, because this enables greater learning. If all sectors immediately adopted and the disaster transpired, many of them would not be able to switch back and avoid the social damages. Gradual adoption allows society to gain from the new technology while updating its beliefs about whether it will have socially damaging uses. As more time passes without disaster, the belief that there will be a disaster declines (“no news is good news”). As society becomes more optimistic, it is optimal to adopt the new technology across a larger number of sectors. Under weak conditions this adoption path is slow and convex, accelerating only after society is fairly certain that a disaster will not occur. A simple quantitative example indicates that, for reasonable parameters for the new technology’s growth advantage and disaster risk, optimal adoption can be very slow.

Perhaps surprisingly, we demonstrate that adoption should be slower when the new technology has a higher growth rate and damages from a disaster are large. This is for two reasons. First, since damages after a potential disaster increase with the new technology’s capabilities, a higher growth rate means that damages also grow more quickly. Second, with a higher growth rate the effective discount rate for future output declines, so that short delays in adoption are not very consequential for discounted utility.

Compared to optimal adoption, equilibrium adoption is inefficiently fast if private firms internalize only part of the social damages from a disaster. Even the order in which sectors adopt the new technology can differ between the equilibrium and the optimum—sectors that have high social damages are not necessarily those that have high *private* damages for adopters.<sup>4</sup>

---

<sup>3</sup>These assumptions can be motivated with generative AI applications. For irreversibility, once large language models like ChatGPT are deployed in secondary education, it may be impossible to roll back their use, even after it becomes clear that they harm student learning. For the damages rising with productivity, many experts fear that these technologies either pose existential risks or will be misused, both of which would be more damaging when they have greater capabilities (e.g., Shevlane, et al., 2023).

<sup>4</sup>For example, if AI is used to create pervasive disinformation on social media, this may be disastrous for

Finally, we discuss how regulatory schemes can help to close the gap between optimal and equilibrium adoption. Pigovian taxes, use taxes, or adoption taxes that are sector-specific can fully implement optimal adoption. When sector-specific policies are not feasible, it is generally not possible to implement optimal technology choices, but regulation can still increase welfare by prohibiting use of the new technology in the sectors with the largest potential for harm until the risk of a disaster is sufficiently low.

This paper is a first attempt to study the consequences and regulation of transformative technologies that can be used for good or bad. Our conclusions naturally depend on our modeling assumptions and should be interpreted with caution.

There are three literatures on which we build. The first is a growing literature on economic disasters (e.g., Rietz, 1988; Barro, 2006, 2009; Weitzman, 2009, 2011; Martin and Pindyck, 2015, 2021), which explores how the risk of rare economic disasters affects asset prices and cost-benefit analysis, but does not focus on questions of technology adoption.

The second is a literature on technology adoption (e.g., Katz and Shapiro, 1986; Parente and Prescott, 1994; Foster and Rosenzweig, 1995, 2010; Acemoglu, Aghion, and Zilibotti, 2006; Acemoglu, Antràs, and Helpman, 2007; Comin and Mestieri, 2014). Early work touching on AI includes Galasso and Luo (2018) and Agrawal, Gans, and Goldfarb (2019), but these papers do not focus on issues of learning about social damages from new technologies.

Third, there is a nascent literature focusing on damages from certain technologies (e.g., Bovenberg and Smulders, 1995; Acemoglu, Aghion, Bursztyn, and Hemous, 2012). Most closely related to our paper are a few works that discuss the dilemma between growth and existential risk from new technologies, including AI. Jones (2023) develops a one-sector growth model in which AI can be used to raise the aggregate growth rate, but with small probability causes human extinction. Whether it is optimal to use AI depends crucially on the coefficient of relative risk aversion and whether consumption utility is bounded. Aschenbrenner (2020) incorporates existential risk into Jones’s (2016) model of growth and mortality, and argues that existential risk rises with consumption unless new mitigation technologies are developed. His model thus exhibits an “existential risk Kuznets curve” in which existential risk optimally increases until sufficient R&D resources are shifted toward mitigation. These two papers share our focus on the costs and benefits of transformative technologies, but they do not address the speed of adoption across sectors and do not feature learning about risks over time.

The rest of the paper is organized as follows. Section I presents our benchmark model. Sections II and III characterize optimal and equilibrium technology choices. Section IV discusses the conditions under which optimal technology choices can be restored through regulatory taxes, and Section V concludes. Omitted proofs and extensions are in the online Appendix.

---

democracy but profitable for social media platforms.

# I Setup

We consider a continuous-time economy that linearly produces a final good from a continuum of sectors  $i \in [0, 1]$ :

$$Y = \int_0^1 Y_i di.$$

A representative household has risk-neutral preferences defined over this final good and discounts the future at rate  $\rho > 0$ .

Each sector can use an old technology  $O$  or a new, transformative technology  $N$ . We write  $Q_j(t) > 0$  for the quality of technology  $j \in \{O, N\}$  at time  $t$ ,  $x_i(t) = 1$  if sector  $i$  switches its production process to technology  $N$ , and  $x_i(t) = 0$  otherwise. Sectoral output is

$$Y_i = (1 - x_i)Q_O + x_i\alpha_iQ_N,$$

where  $\alpha_i$  designates the comparative advantage of the new technology, which may vary if the new technology is better-suited for some sectors than others. Given technology choices  $x = (x_i)_{i \in [0, 1]}$  and qualities  $Q = (Q_O, Q_N)$ , final output is

$$Y(x, Q) = \int_0^1 (1 - x_i)Q_O + x_i\alpha_iQ_N di.$$

The new technology is *transformative*, both because it is general-purpose and can be applied across all sectors, and because it enables not just the production of more output, but a higher growth rate:

$$g_N > g_O \geq 0.$$

As a result of its restructuring impact on the economy, it also poses new risks. We model these by assuming that there may be a *disaster* whereby the new technology generates negative effects. If a disaster happens, then there will be *damages* of  $\delta_i Q_N > 0$  (in units of the final good) in the sectors that are using the technology. We assume that use of the new technology may be irreversible, so that with probability  $\eta_i \in (0, 1)$  sector  $i$  cannot switch to technology  $O$  if it is using technology  $N$  when the disaster strikes. The realization of this reversibility event is independent across sectors. We assume that damages are proportional to  $Q_N$  because the negative effects correspond to misusing the better capabilities of the new technology.

In what follows, we reorder sectors so that  $\delta_i$  is increasing and assume that  $i$  denotes the quantiles of the  $\delta$  distribution, so that we can take this distribution to be uniform over some

interval  $[\underline{\delta}, \bar{\delta}]$ . Overall damages then become

$$D(x, Q) = \left( \int_0^1 \delta_i x_i di \right) Q_N.$$

The economy will experience a disaster with probability  $\bar{\mu} \in (0, 1)$ , and if there is a disaster, its arrival time  $T$  is distributed exponentially with rate  $\lambda$ . We let  $\mu(t)$  denote the (the planner's or society's) posterior belief at  $t$  that there will be a disaster, assuming one has not yet arrived. We impose rational expectations, so that  $\mu(0) = \bar{\mu}$  and the posterior belief evolves according to Bayes's rule:

$$(1) \quad \dot{\mu}(t) = -\lambda\mu(t)(1 - \mu(t)).$$

A few comments are in order. First, we model damages in each sector  $i$  by the reduced-form function  $\delta_i Q_N$  to capture a broad range of potential harms. In the context of AI, these include the spread of disinformation that harms democracy; mass unemployment; and the disruption of production in many sectors from AI-aided cyberattacks.<sup>5</sup> Second, as suggested above, the assumption that damages are proportional to  $Q_N$  is related to the transformative nature of this new technology. For example, damages from disinformation from AI will be higher when it can generate better language. Third, we assume that the arrival rate of the disaster—and hence learning about the negative effects of the new technology—is independent of how many sectors switch to the new technology. This is for simplicity, but is not unreasonable since many of the potential misuses of a new technology can be gradually recognized without widespread adoption.<sup>6</sup> Fourth, it can be verified that our results remain identical if, instead of a single economy-wide disaster, there are sector-specific disasters and beliefs about each sector's disaster follow (1).

## II Socially Optimal Technology Choice

In this section, we set up, solve, and provide comparative statics for the (social) planner's problem.

---

<sup>5</sup>Our functional form assumptions also impose that the rate of substitution between gross consumption and damages in utility is constant and equal to one. Jones (2023) points out that this may not hold in the case of existential risk and explores the implications for optimal use of a life-threatening new technology.

<sup>6</sup>Alternative assumptions are discussed in Section V.

## II.A Social Planner's Problem

Given risk neutrality, the planner's objective is

$$(2) \quad V(0) = \mathbb{E}_{\mu(0)} \left[ \int_0^\infty \exp(-\rho t) [Y(t) - D(t)] dt \right],$$

where  $Y(t)$  and  $D(t)$  denote output and damages at time  $t$  and the expectation  $\mathbb{E}_{\mu(0)}$  is with respect to the prior belief  $\mu(0)$  over the disaster's arrival time  $T$ . To ensure that the objective is well-defined, we assume

$$(3) \quad \rho > g_N,$$

which rules out the case in which the new technology grows so quickly that discounted utility becomes infinite.

It is more convenient to work with the recursive formulation of (2), which has the following state variables: the posterior belief of disaster,  $\mu$ ; the time-varying qualities of the old and new technologies,  $Q$ ; and, after the disaster, the set of sectors that were already using the new technology and for which this use is irreversible. We track these sectors using the vector  $\bar{x} = (\bar{x}_i)_{i \in [0,1]}$ , where  $\bar{x}_i = 1$  if sector  $i$  uses technology  $N$  irreversibly and  $\bar{x}_i = 0$  otherwise. Let  $V(\mu, Q)$  denote pre-disaster social welfare, and let  $W(\bar{x}, Q)$  denote post-disaster welfare. Then the Hamilton-Jacobi-Bellman (HJB) equations for the planner are

$$(4) \quad \rho V(\mu, Q) = \max_{x_i \in \{0,1\}} \{Y(x, Q) + \mu \lambda (\mathbb{E}[W(\bar{x}, Q) | x] - V(\mu, Q))\} + \dot{V}(\mu, Q),$$

$$(5) \quad \rho W(\bar{x}, Q) = \max_{x_i \in \{\bar{x}_i, 1\}} \{Y(x, Q) - D(x, Q)\} + \dot{W}(\bar{x}, Q).$$

Equation (5) imposes that  $x_i$  cannot be less than  $\bar{x}_i$ , because  $\bar{x}_i = 1$  implies that sector  $i$ 's use of the new technology is irreversible.  $V$  then depends on the conditional expectation of welfare after a disaster given the current technology choices  $x$ , denoted by  $\mathbb{E}[W(\bar{x}, Q) | x]$ .<sup>7</sup> In (4) we also use the fact that the arrival rate of the disaster, given the posterior  $\mu$ , is  $\mu \lambda$ .

To characterize the planner's technology choices, suppose first that the disaster has occurred. The planner's problem in (5) is linear, so the solution is

$$x_i = \begin{cases} 1 & \text{if } \bar{x}_i = 1 \text{ or } (\alpha_i - \delta_i)Q_N > Q_O, \\ 0 & \text{else.} \end{cases}$$

---

<sup>7</sup>To determine this conditional expectation, we use  $\mathbb{P}(\bar{x}_i = 1 | x_i = 1) = \eta_i$  and  $\mathbb{P}(\bar{x}_i = 1 | x_i = 0) = 0$ .

This expression assumes, without loss of generality, that the planner sticks with the old technology if indifferent. It also imposes the constraint that  $x_i = 1$  when  $\bar{x}_i = 1$ . Even when unconstrained, it may be optimal to set  $x_i = 1$  if the output produced by technology  $N$  exceeds its damages plus the output that can be produced by technology  $O$ . We first assume that damages are sufficiently large that, whenever possible, the planner chooses technology  $O$  after a disaster:

$$(6) \quad \alpha_i \leq \delta_i.$$

This enables us to focus on the most interesting case where damages exceed the benefits of the new technology. We return to the general case in Section II.C.

Integrating the HJB equation (5) and taking expectations with respect to  $\bar{x}$ , we have

$$\mathbb{E}[W(\bar{x}, Q) | x] = \int_0^1 \left[ (1 - x_i \eta_i) \frac{1}{\rho - g_O} Q_O + x_i \eta_i \frac{\alpha_i - \delta_i}{\rho - g_N} Q_N \right] di.$$

Before the disaster, it is optimal from (4) to use technology  $N$  in sector  $i$  iff

$$(7) \quad \alpha_i Q_N - Q_O > \mu \lambda \eta_i \left[ \frac{1}{\rho - g_O} Q_O - \frac{\alpha_i - \delta_i}{\rho - g_N} Q_N \right].$$

Intuitively, the left-hand side is the flow gain from using technology  $N$  in sector  $i$ , while the right-hand side is the expected loss due to the disaster, including both the discounted value of lost output and the irreversible damages. These losses are multiplied by the posterior arrival rate of the disaster  $\mu \lambda$  and the probability of irreversibility  $\eta_i$ . Since  $\mu$  is decreasing and  $Q_N/Q_O$  is increasing, for any initial state  $(\mu(0), Q(0))$  there exists a time  $t_i < \infty$  such that technology  $O$  is used in sector before  $t_i$  and technology  $N$  is used thereafter.

## II.B Socially Optimal Technology Adoption

To determine how (socially) optimal use of technology  $N$  changes over time, denote the fraction of sectors that use technology  $N$ , or total *adoption*, by

$$X(\mu, q) = \int_0^1 x_i(\mu, q) di.$$

Here  $q = \log(Q_N/Q_O)$  is the *quality gap* between the technologies, and  $x_i(\mu, q) = 1$  iff it is optimal to use technology  $N$  in sector  $i$  in state  $(\mu, q)$ . For simplicity, we assume that  $\alpha_i$  and  $\eta_i$  are constant across sectors and equal to  $\alpha$  and  $\eta$  (the general case is studied in Appendix B).

This implies that there exists a *damage threshold*  $L(\mu, q)$  such that it is optimal to adopt the new technology in sector  $i$  iff  $\delta_i < L(\mu, q)$ . Letting  $F$  denote the cumulative distribution function of the uniform distribution over  $[\underline{\delta}, \bar{\delta}]$ , total adoption is then just the fraction of sectors below the damage threshold:

$$X(\mu, q) = F(L(\mu, q)).$$

The following proposition is immediate from (7), and we omit its proof:

**Proposition 1.** *Suppose (6) holds and  $\alpha_i$  and  $\eta_i$  are constant across sectors. It is socially optimal to use technology  $N$  in sector  $i$  iff  $\delta_i < L(\mu, q)$ , where*

$$(8) \quad \frac{L(\mu, q) - \alpha}{\rho - g_N} = \frac{\alpha - \exp(-q)}{\mu\lambda\eta} - \frac{\exp(-q)}{\rho - g_O}.$$

$L(\mu, q)$  (and thus  $X(\mu, q)$ ) is increasing in  $\alpha$  and  $q$ ; decreasing in  $g_O$ ,  $\lambda$ , and  $\mu$ ; and decreasing in  $g_N$ , provided that  $L(\mu, q) > \alpha$ .

Given (6), the condition  $L(\mu, q) > \alpha$  is satisfied as soon as there is any adoption. Proposition 1 then implies that when the new technology enables *faster growth*, its adoption should be *slower*. This is because of a *precautionary motive*—even though the planner is risk-neutral, she would like to avoid irreversible damages from the new technology. The faster the new technology grows, the greater are the potential net output losses, strengthening this precautionary motive.

The comparative statics in Proposition 1 are partial because they hold the state  $(\mu, q)$  fixed. Full comparative statics must account for how parameter changes affect the evolution of the state  $(\mu(t), q(t))$ . The belief  $\mu(t)$  does not depend on the growth rates  $g_O$  and  $g_N$ , but the quality gap  $q(t) = q(0) + (g_N - g_O)t$  does. The damage threshold  $L(\mu, q)$  is increasing in the quality gap, so any change in growth rates affects adoption at each  $t > 0$  through both the direct effects described in Proposition 1 and the indirect effects through changes in the quality gap  $q(t)$ . The next proposition characterizes these total effects.

**Proposition 2.** *Suppose (6) holds and  $\alpha_i$  and  $\eta_i$  are constant across sectors.*

1.  $X(\mu(t), q(t))$  is decreasing in  $g_O$ .
2. There exists an earliest time  $\bar{t} < \infty$  such that  $X(\mu(t), q(t))$  is decreasing in  $g_N$  if  $t > \bar{t}$ . The time  $\bar{t}$  is decreasing in  $g_N$ .
3. Adoption falls to zero as  $g_N$  approaches  $\rho$ , i.e.,  $\lim_{g_N \uparrow \rho} X(\mu(t), q(t)) = 0$ .



The first part of Proposition 2 establishes that the comparative static for  $g_O$  from Proposition 1 generalizes in the presence of the indirect effects through  $q(t)$ —the quality gap  $q(t)$  is declining in  $g_O$ , reinforcing the direct effect and decreasing adoption. The second part shows that the new technology's growth rate has more nuanced implications: Adoption is not always decreasing in  $g_N$ , but it is after some critical time  $\bar{t}$ , and this time itself is a decreasing function of  $g_N$ . This holds because the precautionary motive highlighted above must compete with the fact that the quality gap  $q(t)$  is increasing in  $g_N$ , but this indirect effect can dominate only at short time horizons.

The third part of proposition establishes that as  $g_N$  increases towards the discount rate, adoption almost stops. This might appear paradoxical initially, but is also intuitive. When  $g_N$  is approximately equal to  $\rho$ , the benefits from the new technology are very high, leading to nearly infinite discounted utility provided no disaster arrives. Delay in adoption thus has little effect on these benefits. However, a disaster will have huge negative consequences, and avoiding it now takes precedence.

The next proposition further characterizes the shape of the adoption curve. Since  $F$  is uniform,  $\dot{X}(\mu, q) = f \dot{L}(\mu, q)$ , where  $f$  is the constant density of  $F$ . Hence, the *curvature* of technology adoption is

$$\frac{\ddot{X}(\mu, q)}{\dot{X}(\mu, q)} = \frac{\ddot{L}(\mu, q)}{\dot{L}(\mu, q)}.$$

We therefore have:

**Proposition 3.** *Suppose (6) holds.*

1.  $\dot{L}(\mu, q) > 0$  is decreasing in  $g_O$ , and it is decreasing in  $g_N$  iff the quality gap is sufficiently large, i.e.,

$$\alpha \exp(q) - 1 > \frac{(\rho - g_N) - (g_N - g_O)}{1 - \mu} \left( \frac{1}{\lambda} + \frac{\mu\eta}{\rho - g_O} \right).$$

2. There exists a positive constant  $G(\mu, q)$  such that if  $\alpha \exp(q) > 1$ ,  $\ddot{L}(\mu, q) > 0$  is positive iff  $g_N - g_O > G(\mu, q)$ .  $G(\mu, q)$  is independent of  $g_N$  and increases to infinity over time.

The intuition for the first part is the same as for Proposition 2: The damage threshold increases as the posterior belief  $\mu$  falls and the quality gap  $q$  grows. Faster growth for technology  $O$  slows the rate of increase of the quality gap and raises the opportunity cost of using technology  $N$  after the disaster. Consequently, the damage threshold grows less quickly in each state. Faster growth for technology  $N$  raises both the rate of increase in the quality gap and the net output losses from technology  $N$  after the disaster. The latter effect dominates when

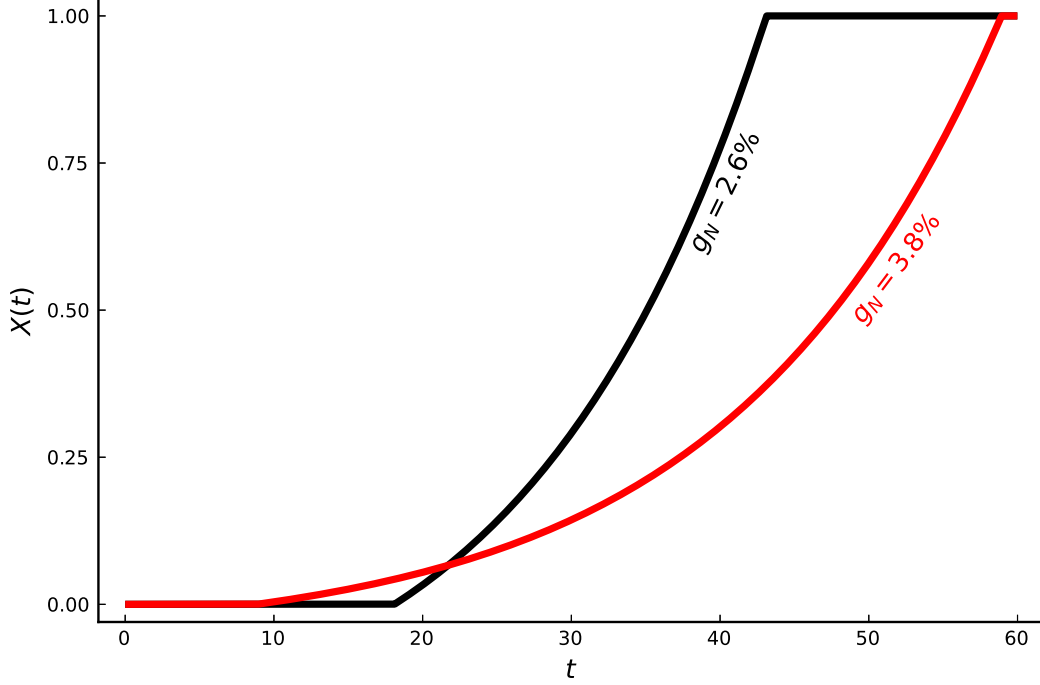


Figure 1: Adoption curves  $X(t) \equiv X(\mu(t), q(t))$  for different values of  $g_N$ . The remaining parameter values are  $\rho = 0.04$ ,  $\lambda = 0.05$ ,  $\eta = 0.5$ ,  $\alpha = 1$ ,  $g_O = 0.02$ ,  $\underline{\delta} = 1$ , and  $\bar{\delta} = 5$ . The initial state is  $\mu(0) = 0.2$  and  $q(0) = 0$ .

the quality gap is sufficiently large because additional improvements in technology  $N$  relative to  $O$  have only a negligible impact on the planner's technology choice.<sup>8</sup>

The second part of the proposition proves that when the new technology's growth advantage is sufficiently large, its adoption will have a convex segment where adoption accelerates. This result holds even though the learning rate  $|\dot{\mu}|$  falls at a greater than exponential rate when  $\mu < \frac{1}{2}$  (in particular,  $\frac{d}{dt} |\dot{\mu}| = -\lambda |\dot{\mu}| (1 - 2\mu)$ ). This is because expected damages from technology  $N$  in sector  $i$  are proportional to the posterior  $\mu$ , and as  $\mu$  declines, larger increases in the damage threshold  $L(\mu, q)$  are needed to balance the expected damages and benefits in the “marginal” sector.<sup>9</sup>

To illustrate these results, we depict the time path of adoption in a couple of parameterized cases in Figure 1. We set  $g_O = 2\%$  in line with trend GDP growth in developed economies and  $\rho = 0.04$  to produce a risk-free interest rate of 4%. We choose two values for  $g_N$  based on Chui, Roberts, Yee, Hazan, et al. (2023), who forecast an increase in the growth rate of 0.6-3.6% in the United States between 2023 and 2040 from AI and other automation technologies. We take the lower end of this range,  $g_N - g_O = 0.6\%$ , and a higher but still conservative estimate

<sup>8</sup>The latter effect also dominates regardless of the quality gap whenever  $L(\mu, q) > 0$  and  $g_N - g_O \geq \rho - g_N$ .

<sup>9</sup>In Appendix B, we verify this intuition by showing that learning dynamics favor *concave* adoption when sectors are heterogeneous according to  $\alpha_i$  instead of  $\delta_i$ .

from the middle of the range:  $g_N - g_O = 1.8\%$  (while still satisfying (3)). We take the two technologies to have the same quality in year  $t = 0$ , thus  $q(0) = 0$ . We suppose that damages range from one to five times gross sectoral output ( $\underline{\delta} = 1, \bar{\delta} = 5$ ), and we set  $\eta = 0.5$  so that half of all sectors using the new technology cannot switch back after a disaster. We set the expected arrival time of a disaster (if one exists) to be 20 years, which gives  $\lambda = 0.05$ . Finally, a recent survey of AI experts reports a median estimate of existential risk of about 10%,<sup>10</sup> and since we are interested in non-existential misuses of AI as well, we choose the initial disaster probability to be twice as large,  $\mu(0) = 20\%$ . Figure 1 shows that optimal adoption is slow, taking about 40 years until full adoption when  $g_N = 2.6\%$  and almost 60 years when  $g_N = 3.8\%$ .

## II.C Optimal Adoption with Small Damages

We have so far imposed (6), ensuring that the post-disaster damages from the new technology are *large* and exceed its gross output within each sector. This is a natural benchmark, since our analysis is motivated by significant potential harms from AI. We now relax this assumption and allow a sector's damages to be *small* relative to its output under the new technology ( $\delta_i < \alpha$ ).

In Appendix C, we show that socially optimal adoption is again characterized by a damage threshold  $L(\mu, q)$ , and we prove the following analogue to Proposition 2 for small damages.

**Proposition 4.** *Suppose  $\alpha_i$  and  $\eta_i$  are constant across sectors. For all  $t$  with  $L(\mu(t), q(t)) < \alpha$ :*

1.  $X(\mu(t), q(t))$  is decreasing in  $g_O$ .
2.  $X(\mu(t), q(t))$  is increasing in  $g_N$ .
3. If  $q(0)$  is sufficiently low and  $X(\mu(t), q(t)) < F(\alpha)$ , adoption is bounded below  $F(\alpha)$  as  $g_N$  approaches  $\rho$ , i.e.,  $\lim_{g_N \uparrow \rho} X(\mu(t), q(t)) < F(\alpha)$ .

Adoption among sectors with small damages is still decreasing in  $g_O$ , but in contrast to the case with large damages, it is increasing in  $g_N$ . Gradual adoption remains optimal even when  $g_N$  increases toward the discount rate  $\rho$ . With small damages, using technology  $N$  is always optimal in the long run. Nevertheless, gradual adoption is optimal to learn about the probability of a disaster (before one occurs) and to delay the adoption of technology  $N$  in case of a disaster until the quality gap becomes sufficiently large. This strategy thus avoids *temporary* costs of irreversibility. Further analysis of this case is presented in Appendix C.

Finally, we note that if damages are uncertain, any chance of large damages leads to longer optimal delay, even if expected damages are small, in order to avoid the possibility that damages turn out to be large *and* adoption is irreversible.

---

<sup>10</sup><https://aiimpacts.org/2022-expert-survey-on-progress-in-ai>

In summary, the optimal adoption of a new, transformative technology should be gradual, particularly when its superior capabilities also make its potential damages greater and there is learning about the likelihood of misuse (a “disaster”).

### III Equilibrium Technology Choice

We now characterize equilibrium technology adoption when private firms do not fully internalize social damages.

#### III.A The Firm’s Problem

Suppose now that in each sector, the choice of technology is made by a private (representative) firm that seeks to maximize expected discounted profits. To simplify, we assume that the firm in sector  $i$  appropriates all output of its intermediate as profits, but only internalizes *private damages*  $\gamma_i \leq \delta_i$ . This textbook externality leads to excessively fast adoption of the new technology before the disaster, and our main results below describe how the equilibrium and socially optimal adoption curves differ.

Firm  $i$ ’s profit maximization problem can be formulated recursively in the same way as the planner’s problem in the previous section. The state variables before the disaster are again  $\mu$  and  $Q$ , and after the disaster they are  $\bar{x}_i$  and  $Q$ . Let  $\Pi_i(\mu, Q)$  denote the firm’s pre-disaster value,  $\Phi_i(\bar{x}_i, Q)$  its post-disaster value, and  $Y_i(x_i, Q)$  its (gross) output. The HJB equations for the firm are

$$(9) \quad \rho \Pi_i(\mu, Q) = \max_{x_i \in \{0,1\}} \{Y_i(x_i, Q) + \mu \lambda (\mathbb{E}[\Phi_i(\bar{x}_i, Q) | x_i] - \Pi_i(\mu, Q))\} + \dot{\Pi}_i(\mu, Q),$$

$$(10) \quad \rho \Phi_i(\bar{x}_i, Q) = \max_{x_i \in \{\bar{x}_i, 1\}} \{Y_i(x_i, Q) - x_i \gamma_i Q_N\} + \dot{\Phi}_i(\bar{x}_i, Q).$$

These value functions differ from the planner’s (4) and (5) because the firm internalizes only a fraction  $\gamma_i/\delta_i$  of the flow damages from technology  $N$ .

We now impose a stronger version of (6): private damages are sufficiently large that firm  $i$  will always choose technology  $O$  after the disaster if possible:<sup>11</sup>

$$(11) \quad \alpha_i \leq \gamma_i.$$

---

<sup>11</sup>Without this assumption, an additional inefficiency would arise in equilibrium as firms would use the new technology in some (reversible) sectors even after a disaster.

Similar to the planner's solution, it is privately optimal for firm  $i$  to use technology  $N$  iff

$$\alpha_i Q_N - Q_O > \mu \lambda \eta_i \left[ \frac{1}{\rho - g_O} Q_O - \frac{\alpha_i - \gamma_i}{\rho - g_N} Q_N \right].$$

The only difference between this condition and the planner's optimality condition (7) is that private damages  $\gamma_i$  appear instead of social damages  $\delta_i$  on the right-hand side. Firm  $i$  internalizes fewer damages from technology  $N$  and thus begins using it earlier.

### III.B Equilibrium Technology Adoption

We denote total equilibrium adoption by

$$\tilde{X}(\mu, q) = \int_0^1 \tilde{x}_i(\mu, q) di,$$

where  $\tilde{x}_i(\mu, q) = 1$  iff firm  $i$  uses technology  $N$  in state  $(\mu, q)$ . Again assuming that  $\alpha_i$  and  $\eta_i$  are constant across sectors, it is immediate that firm  $i$  will adopt the new technology iff *private* damages are lower than the damage threshold,  $\gamma_i < L(\mu, q)$ . Equilibrium adoption is then

$$\tilde{X}(\mu, q) = F_\gamma(L(\mu, q)),$$

where  $F_\gamma$  is the cumulative density function of  $\gamma_i$ .

This characterization implies that all comparative statics results from Section II.B apply to equilibrium adoption. The results in Propositions 1 and 3 concern only the damage threshold  $L(\mu, q)$  and hold exactly as stated, while Proposition 2 applies after replacing  $X(\mu, q)$  with  $\tilde{X}(\mu, q)$ :

**Proposition 5.** *Suppose (11) holds and  $\alpha_i$  and  $\eta_i$  are constant across sectors.*

1.  $\tilde{X}(\mu(t), q(t))$  is decreasing in  $g_O$ .
2. There exists an earliest time  $\tilde{t} < \infty$  such that  $\tilde{X}(\mu(t), q(t))$  is decreasing in  $g_N$  if  $t > \tilde{t}$ . The time  $\tilde{t}$  is decreasing in  $g_N$ .
3. Adoption falls to zero as  $g_N$  increases to  $\rho$ :  $\lim_{g_N \uparrow \rho} \tilde{X}(\mu(t), q(t)) = 0$ .

In the remainder of this section, we characterize how the optimal and equilibrium adoption curves differ. We first observe that similar adoption curves do not imply that the equilibrium is optimal, because the order in which sectors adopt the new technology matters. For example,

private and social damages may be *negatively affiliated*, meaning that high social damage sectors have low private damages. In this case, the order in which the new technology spreads in equilibrium is exactly the opposite of the optimal order.

Even when the equilibrium and optimal orders of adoption coincide, the equilibrium can be inefficient. To see this, suppose that social and private damages are *positively affiliated*, so that there exists a non-negative and (strictly) increasing function  $\kappa$  with  $\gamma_i = \kappa(\delta_i) \leq \delta_i$ . We can then write equilibrium adoption as

$$\tilde{X}(\mu, q) = F(\kappa^{-1}(L(\mu, q))).$$

This equation implies that the equilibrium adoption curve  $\tilde{X}(\mu(t), q(t))$  is a distorted version of the optimal adoption curve, with an *equilibrium damage threshold*  $\tilde{L}(\mu, q) = \kappa^{-1}(L(\mu, q))$ . In this case, knowing how the equilibrium and social damage thresholds differ is sufficient to fully characterize equilibrium inefficiencies. The next proposition determines how the level, rate of change, and curvature of the equilibrium damage threshold  $\tilde{L}(\mu, q)$  differ from its social counterpart  $L(\mu, q)$ .

**Proposition 6.** *Suppose (11) holds and  $\alpha_i$  and  $\eta_i$  are constant across sectors.*

1. *The equilibrium damage threshold is always greater than the social damage threshold:*  
 $\tilde{L}(\mu, q) \geq L(\mu, q)$ .
2. *The equilibrium damage threshold increases more quickly than the social damage threshold when  $\kappa'(\tilde{L}(\mu, q)) < 1$ :*

$$\dot{\tilde{L}}(\mu, q) = \frac{\dot{L}(\mu, q)}{\kappa'(\tilde{L}(\mu, q))}.$$

3. *The equilibrium damage threshold is more convex than the social damage threshold when  $\kappa$  is locally concave:*

$$\frac{\ddot{\tilde{L}}(\mu, q)}{\dot{\tilde{L}}(\mu, q)} = \frac{\ddot{L}(\mu, q)}{\dot{L}(\mu, q)} - \frac{\kappa''(\tilde{L}(\mu, q))}{\kappa'(\tilde{L}(\mu, q))} \dot{L}(\mu, q).$$

These results follow from the definition of the equilibrium damage threshold  $\tilde{L}(\mu, q)$ . We illustrate them in Figure 2 by depicting socially optimal and equilibrium adoption curves for the benchmark parameterizations in Figure 1 and a concave affiliation function  $\kappa$ . The equilibrium damage threshold is always greater than its social counterpart and increases more quickly (for the marginal sectors where  $\kappa'(\tilde{L}(\mu, q)) < 1$ ). Consequently, equilibrium adoption is inefficiently rapid and accelerates when there are high social damages.

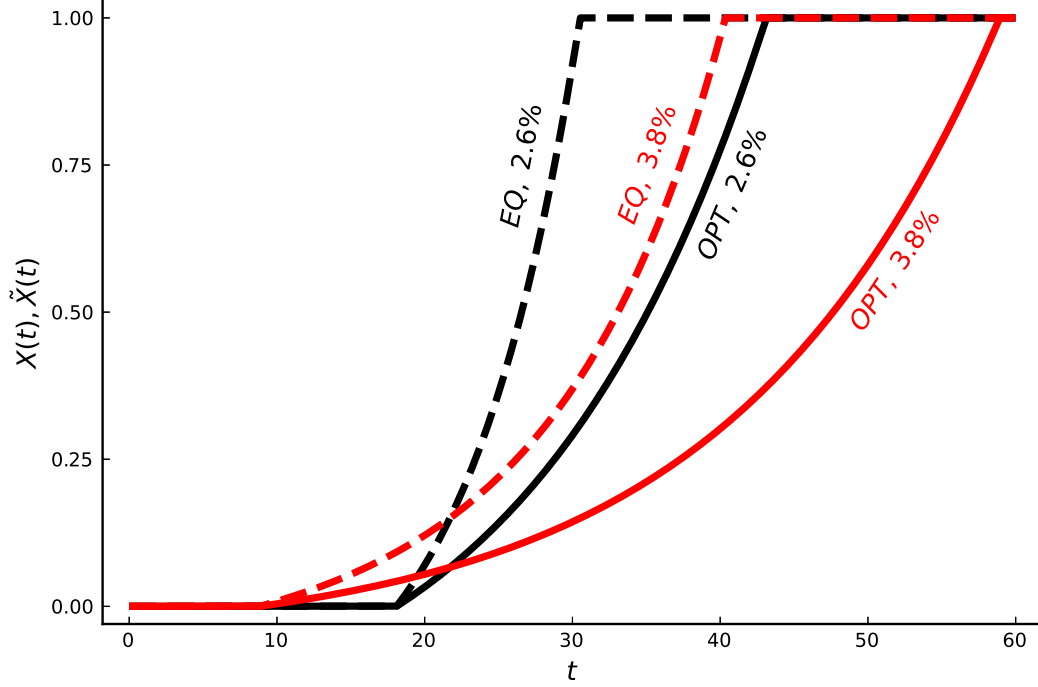


Figure 2: Socially optimal and equilibrium adoption curves,  $X(t)$  and  $\tilde{X}(t) \equiv \tilde{X}(\mu(t), q(t))$ . The calibration is the same as in Figure 1. The affiliation function is  $\kappa(\delta) = \delta^{1/2}$ .

In summary, equilibrium adoption of transformative technologies is determined by the same forces that shape optimal adoption. However, because firms are motivated by higher productivity and discouraged only by private damages, equilibrium adoption is generally suboptimal: Firms do not fully internalize social damages from potential disasters, so equilibrium adoption is typically too high and rises too quickly, and the order in which sectors adopt the new technology may differ from the optimal one.

## IV Regulating Technology Choice

Since equilibrium adoption is potentially inefficient, a natural question is whether government regulation can close the gap between equilibrium and optimal adoption decisions. Throughout this section, we continue to assume that (11) holds, and we simplify the analysis by focusing on *ex ante* regulations.<sup>12</sup>

Socially optimal and equilibrium technology choices differ because the planner and private firms internalize different damages after the disaster and hence different *expected* damages

<sup>12</sup>We ignore *ex post* (“Pigovian”) taxes both because their analysis is essentially identical to our characterization of use taxes, and also because they may not be credible as they do not affect technology choice after the disaster—the private sector already stops using the new technology whenever possible.

before the disaster. A straightforward way to correct firms' incentives is through a *use tax* that raises firms' costs of using the new technology *before the disaster*.<sup>13</sup> When sector-specific taxes are feasible, the tax that implements the optimal technology choice for sector  $i$  is equal to the difference between expected discounted social and private damages:

$$(12) \quad \tau_i(\mu, Q_N) = \mu \lambda \eta_i \frac{\delta_i - \gamma_i}{\rho - g_N} Q_N.$$

The next proposition notes several properties of these optimal taxes.

**Proposition 7.** *The optimal use tax  $\tau_i(\mu, Q_N)$  is larger in sectors with a larger probability of irreversibility  $\eta_i$  and a larger difference between social and private damages  $\delta_i - \gamma_i$ . It is log-concave in time and limits to zero as  $t \rightarrow \infty$  iff  $\lambda > g_N$ .*

The cross-sector comparative statics follow immediately from (12). Differentiating (12) with respect to time yields

$$\frac{\dot{\tau}_i(\mu, Q_N)}{\tau_i(\mu, Q_N)} = \frac{\dot{\mu}}{\mu} + \frac{\dot{Q}_N}{Q_N} = -\lambda(1 - \mu) + g_N.$$

Since  $\mu$  declines before the disaster,  $\tau_i(\mu(t), Q_N(t))$  is log-concave. The difference between social and private damages from a disaster is increasing in  $Q_N$ , pushing taxes higher, while growing optimism about the absence of a disaster pushes taxes lower. The tax is eventually decreasing to zero iff learning about the disaster risk is sufficiently fast,  $\lambda > g_N$ .

Sector-specific taxes require detailed information about damages and may generally be difficult to implement. Even in the benchmark case in which  $\alpha_i$  and  $\eta_i$  are constant across sectors, the next proposition shows that a sector-independent tax scheme cannot correct inefficient equilibrium adoption unless social and private damages are positively affiliated.

**Proposition 8.** *Suppose  $\alpha_i$  and  $\eta_i$  are constant across sectors. Given any sector-independent use tax  $\tau(\mu, Q)$ , firm  $i$  begins using technology  $N$  earlier than firm  $j$  iff  $\gamma_i \leq \gamma_j$ . Socially optimal technology choices can be implemented for any initial state  $(\mu(0), Q(0))$  iff social and private damages are positively affiliated. In this case, the following tax is optimal:*

$$(13) \quad \tau(\mu, Q) = \mu \lambda \eta \frac{L(\mu, q) - \kappa(L(\mu, q))}{\rho - g_N} Q_N.$$

This proposition clarifies that a sector-independent tax can differentially delay adoption for sectors with different private damages  $\gamma_i$ , but it cannot alter the order of adoption. When

---

<sup>13</sup>Naturally, *adoption taxes* that are paid when new technologies are first introduced are equivalent.



private and social damages are positively affiliated, the socially optimal and equilibrium orders of adoption coincide, so a sector-independent tax can fully correct equilibrium inefficiencies.

When the optimal and equilibrium orders of adoption differ, a different policy that we refer to as a *regulatory sandbox* may be more effective. Under this policy, sectors with social damages below a threshold  $\hat{\delta}$  (“inside the sandbox”) can choose their technology freely, while sectors above the threshold are restricted from using the new technology until time  $\hat{T}$ . This policy allows the planner to ensure that sectors with high social damages adopt only after the new technology is established to be relatively safe. The next proposition demonstrates that the sandbox policy can improve upon the laissez-faire equilibrium.

**Proposition 9.** *Suppose  $\alpha_i$  and  $\eta_i$  are constant and  $\gamma_i < \delta_i$  across sectors. Then there exists a sandbox policy  $(\hat{\delta}, \hat{T})$  that strictly improves upon the laissez-faire equilibrium.*

In Appendix D, we provide additional details about optimal regulatory sandboxes and compare them to sector-independent taxes. In general, each of these policies can improve upon the laissez-faire equilibrium, and combining both is better: A sector-independent tax can differentially delay adoption for sectors with varying private damages  $\gamma_i$ , but it cannot alter the order of adoption. A regulatory sandbox can alter the order by delaying adoption for sectors with high social damages.

To implement welfare-improving use taxes or regulatory sandboxes, regulators must have some knowledge about the potential social damages from the new technology across different sectors. Although there is still substantial uncertainty about these damages from AI, assuming some knowledge is reasonable and consistent with current approaches to regulation. For example, the EU AI Act proposal outlines a framework in which an AI system is subject to different regulations depending on whether its intended use is considered “high risk” (e.g., the operation of critical infrastructure, employment and worker management processes, or law enforcement; see European Commission, 2021). The results in this section provide a foundation for this regulatory approach and also suggest that these policies should be updated as AI technologies become more capable and as society learns more about the risks.

## V Concluding Remarks

Advances in generative AI technologies, such as large language models, have intensified both hopes of more rapid economic growth and concerns about their potential negative consequences. Despite a robust public discussion on AI, there are currently no economic models of the regulation of transformative technologies. This paper has taken a first step in building such a model to provide novel insights for this debate.

We consider the adoption decision of a new, transformative technology that can increase productivity growth across all sectors of the economy but also raises risks of misuse, which we model as the stochastic arrival of a “disaster”. If a disaster occurs, some of the sectors using the new technology may be unable to switch back to the old, safe technology. Whether a disaster will occur is unknown, and society gradually learns about it over time. Consequently, adoption should be gradual and typically follows a convex path, initially growing slowly before accelerating later. Most surprisingly, a faster growth rate of the new technology should lead to slower adoption when potential damages are large: Although the planner is risk-neutral, she has a *precautionary motive* as irreversible damages imply that it is better to wait and learn about the likelihood of a disaster. These irreversible damages are greater when the new technology has a higher growth rate, strengthening the precautionary motive. Finally, if private firms internalize only part of the social damages from transformative technologies, equilibrium adoption is too fast and necessitates regulatory policies.

There are many interesting areas left for future work. First, in contrast to our baseline assumptions, early adoption may increase risks or may facilitate either general learning about potential misuses of the new technology or sector-specific learning about its “safe use”. These considerations may motivate “experimentation” by adopting the technology in a few sectors or trying different uses in some areas, which is an important topic for future work.

Second, many of the misuses of new AI technologies depend on market structure and other aspects of regulation (e.g., concerning disinformation, discrimination, or privacy), and it would be interesting to explore how these affect optimal and equilibrium adoption.

Third, we simplified the analysis by assuming risk neutrality. Jones (2023) demonstrates that the extent of risk aversion and the precise form of damages have a first-order effect on the trade-off between higher growth and the likelihood of a disaster, and these can be incorporated in future analyses of learning about misuses of new transformative technologies.

Fourth, we abstracted from choices about how new technologies may be used. If regulations or other factors can prevent misuse of technology, then faster adoption can become optimal.

Finally, we showed that the optimal path of adoption depends on a few parameters, but there is currently a huge amount of uncertainty about their values. Careful empirical assessment of the costs and benefits of new, transformative technologies like generative AI is an obvious area for future research.

## References

- Acemoglu, D., Aghion, P., Bursztyn, L., & Hémous, D. (2012). The environment and directed technical change. *American economic review*, 102(1), 131–166.
- Acemoglu, D., Aghion, P., & Zilibotti, F. (2006). Distance to frontier, selection, and economic growth. *Journal of the European Economic association*, 4(1), 37–74.
- Acemoglu, D., Antràs, P., & Helpman, E. (2007). Contracts and technology adoption. *American Economic Review*, 97(3), 916–943.
- Agrawal, A., Gans, J., & Goldfarb, A. (2019). Economic policy for artificial intelligence. *Innovation policy and the economy*, 19(1), 139–159.
- Aschenbrenner, L. (2020). *Existential risk and growth* (tech. rep.). GPI Working Paper.
- Barro, R. J. (2006). Rare disasters and asset markets in the twentieth century. *The Quarterly Journal of Economics*, 121(3), 823–866.
- Barro, R. J. (2009). Rare disasters, asset prices, and welfare costs. *American Economic Review*, 99(1), 243–264.
- Bovenberg, A. L., & Smulders, S. (1995). Environmental quality and pollution-augmenting technological change in a two-sector endogenous growth model. *Journal of public Economics*, 57(3), 369–391.
- Chui, M., Roberts, R., Yee, L., Hazan, E., Singla, A., Smaje, K., Sukharevsky, A., & Zimmel, R. (2023). *The economic potential of generative ai* (tech. rep.). McKinsey & Company.
- Comin, D., & Mestieri, M. (2014). Technology diffusion: Measurement, causes, and consequences. In *Handbook of economic growth* (pp. 565–622). Elsevier.
- European Commission. (2021). Proposal for a regulation of the european parliament and of the council laying down harmonised rules on artificial intelligence (artificial intelligence act) and amending certain union legislative acts. <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=celex:52021PC0206>
- Foster, A. D., & Rosenzweig, M. R. (1995). Learning by doing and learning from others: Human capital and technical change in agriculture. *Journal of political Economy*, 103(6), 1176–1209.
- Foster, A. D., & Rosenzweig, M. R. (2010). Microeconomics of technology adoption. *Annu. Rev. Econ.*, 2(1), 395–424.
- Galasso, A., & Luo, H. (2018). Punishing robots: Issues in the economics of tort liability and innovation in artificial intelligence. In *The economics of artificial intelligence: An agenda* (pp. 493–504). University of Chicago Press.
- Jones, C. I. (2016). Life and growth. *Journal of political Economy*, 124(2), 539–578.

- Jones, C. I. (2023). *The ai dilemma: Growth versus existential risk* (tech. rep.). Stanford GSB. Mimeo.
- Katz, M. L., & Shapiro, C. (1986). Technology adoption in the presence of network externalities. *Journal of political economy*, 94(4), 822–841.
- Martin, I. W., & Pindyck, R. S. (2015). Averting catastrophes: The strange economics of scylla and charybdis. *American Economic Review*, 105(10), 2947–2985.
- Martin, I. W., & Pindyck, R. S. (2021). Welfare costs of catastrophes: Lost consumption and lost lives. *The Economic Journal*, 131(634), 946–969.
- Parente, S. L., & Prescott, E. C. (1994). Barriers to technology adoption and development. *Journal of political Economy*, 102(2), 298–321.
- Rietz, T. A. (1988). The equity risk premium a solution. *Journal of monetary Economics*, 22(1), 117–131.
- Shevlane, T., Farquhar, S., Garfinkel, B., Phuong, M., Whittlestone, J., Leung, J., Kokotajlo, D., Marchal, N., Anderljung, M., Kolt, N., Ho, L., Siddarth, D., Avin, S., Hawkins, W., Kim, B., Gabriel, I., Bolina, V., Clark, J., Bengio, Y., ... Dafoe, A. (2023). Model evaluation for extreme risks.
- Weitzman, M. L. (2009). On modeling and interpreting the economics of catastrophic climate change. *The review of economics and statistics*, 91(1), 1–19.
- Weitzman, M. L. (2011). Fat-tailed uncertainty in the economics of catastrophic climate change. *Review of Environmental Economics and Policy*.