

# The Existential Risks of Artificial Intelligence: Analyzing p(doom) and Recent Developments

The rapid advancement of artificial intelligence (AI) has sparked significant debate and concern among researchers, policymakers, and the general public. One of the most pressing topics in this discourse is the concept of "p(doom)," which stands for the probability of catastrophic outcomes, or "doom," resulting from AI. This term, originating as an inside joke among AI researchers, has gained prominence following the release of advanced AI models like GPT-4. High-profile figures such as Geoffrey Hinton and Yoshua Bengio have publicly warned about the existential risks posed by AI, contributing to the growing anxiety surrounding this technology.

Recent reports and articles have highlighted various scenarios in which AI could potentially lead to human extinction or severe societal disruption. For instance, a [report commissioned by the US State Department](#) detailed the catastrophic risks associated with artificial general intelligence (AGI), a yet-to-be-achieved level of AI that could outperform humans across a broad range of domains. The report emphasized the potential for AI weaponization, including biowarfare, mass cyber-attacks, and autonomous robots, which could pose an existential threat to humanity.

Moreover, a [survey conducted by AI Impacts](#) revealed that a significant portion of AI researchers believe there is a non-negligible risk of AI causing human extinction. The survey's median participant estimated a 5% likelihood of AI-driven extinction by the year 2100, reflecting the widespread concern within the AI research community.

Prominent AI experts and industry leaders have also voiced their concerns. For example, [Sam Altman, CEO of OpenAI](#), has expressed fears about the potential for AI to lead to catastrophic outcomes, while [Elon Musk](#) has predicted a 10-20% chance that AI will destroy humanity. These predictions underscore the urgency of addressing the risks associated with AI development.

In response to these concerns, various organizations and experts have called for increased regulation and safety measures. The [Center for AI Safety](#) released a

statement signed by nearly 400 AI experts, emphasizing that mitigating the risk of AI-induced extinction should be a global priority alongside other societal-scale risks such as pandemics and nuclear war.

As the debate over AI's existential risks continues, it is crucial to critically examine the potential consequences of advanced AI systems and develop robust strategies to mitigate these risks. This report aims to provide a comprehensive analysis of the concept of p(doom), review recent developments and expert opinions, and explore the implications for public policy and future research.

## Table of Contents

- Understanding p(doom) and Its Origins
  - Definition and Context of p(doom)
  - Expert Opinions and Probability Estimates
  - Factors Contributing to p(doom)
  - Criticisms and Debates
  - Case Studies and Real-World Examples
  - Public Perception and Cultural Impact
  - Survey Data and Research Findings
  - Mitigation Strategies and Future Directions
  - Conclusion
- Recent Reports and Expert Opinions on AI Risks
  - Expert Opinions on AI Risks
    - Probability Estimates and p(doom)
    - Recent Surveys and Reports
  - Potential Catastrophic Scenarios
    - Biological Warfare and Cyber Attacks
    - Autonomous Weapons and Misinformation
  - Mitigation Strategies and Policy Recommendations
    - Regulation and Oversight
    - Research and Community Resilience
  - Diverging Opinions and Debates
    - Skepticism and Optimism
    - Incremental Harms and Long-Term Risks
- Potential Scenarios and Mitigation Strategies
  - Potential Scenarios of AI-Induced Catastrophes
    - Full Automation of Labor and Economic Disruption

- Autonomous Weapons and Military Conflicts
- Cybersecurity Threats and Infrastructure Attacks
- Disinformation and Social Manipulation
- Mitigation Strategies for AI Risks
  - Regulation and Oversight
  - Research and Development of Safe AI
  - Scenario Planning and Risk Assessment
  - International Collaboration and Governance
  - Public Awareness and Education
- Conclusion

## Understanding p(doom) and Its Origins

### Definition and Context of p(doom)

The term "p(doom)" refers to the probability of catastrophic outcomes, or "doom," resulting from artificial intelligence (AI). This concept is particularly concerned with the existential risks posed by artificial general intelligence (AGI), which is AI that can perform any intellectual task that a human can ([Wikipedia](#)). The term originated as an inside joke among AI researchers but gained prominence in 2023 following the release of GPT-4. High-profile figures such as Geoffrey Hinton and Yoshua Bengio began to publicly warn about the risks associated with AI ([ABC News](#)).

### Expert Opinions and Probability Estimates

Several AI experts have provided their estimates for p(doom). For instance, Yoshua Bengio, a prominent AI researcher, has estimated a 20% probability that AI could lead to catastrophic outcomes. His estimate is based on a 50% probability that AI will reach human-level capabilities within a decade and a greater than 50% likelihood that AI or humans will misuse the technology ([ABC News](#)). Similarly, Elon Musk has suggested a 10-20% chance that AI might end humanity, although he believes the potential positive outcomes of AI development outweigh the negative scenarios ([Windows Central](#)).

On the more extreme end, Roman Yampolskiy, an AI safety researcher, has argued that the probability of AI ending humanity is much higher, estimating it at 99.999999% ([Windows Central](#)).

## Factors Contributing to p(doom)

Several factors contribute to the varying estimates of p(doom). These include the speed of AI development, the potential for AI to reach or exceed human-level intelligence, and the likelihood of AI being used maliciously or uncontrollably. For example, the rapid advancements in generative AI and large language models like OpenAI's ChatGPT have heightened concerns about AI's potential risks ([Futurism](#)).

Additionally, geopolitical factors play a role. The US government's export rules preventing chipmakers like NVIDIA from shipping chips to China are partly motivated by concerns about AI being used in military advances ([Windows Central](#)).

## Criticisms and Debates

The concept of p(doom) has faced criticism and sparked debate among experts. One major point of contention is the lack of clarity about whether predictions are conditional on the existence of AGI, the time frame considered, and the precise definition of "doom" ([Wikipedia](#)). Critics argue that the term can be overly vague and alarmist, potentially overshadowing more immediate and tangible issues related to AI, such as privacy concerns and ethical considerations ([Futurism](#)).

## Case Studies and Real-World Examples

To better understand the implications of p(doom), it is useful to examine real-world case studies. For instance, Anthropic, an AI research organization, focuses on building reliable, interpretable, and steerable AI systems. Their comprehensive strategy aims to mitigate AI risks and significantly reduce the probability of catastrophic outcomes ([Podtail](#)).

Another example is the case of Sony's AIBO robot dog, which simulates emotional responses. This case study highlights the ethical considerations and potential risks associated with AI emotionality and the future of intelligent machines ([Podtail](#)).

## Public Perception and Cultural Impact

The concept of p(doom) has permeated popular culture and public discourse. It is often discussed in media and among tech enthusiasts, sometimes even becoming a topic of casual conversation. For example, Aaron Levie, CEO of the cloud platform

Box, mentioned that  $p(\text{doom})$  comes up in almost every dinner conversation among tech professionals ([Futurism](#)).

Moreover, the term has been popularized by forums like LessWrong, where it was first used by a programmer named Tim Tyler to refer to the probability of doom without being too specific about the time scale or the definition of "doom" ([Futurism](#)).

## Survey Data and Research Findings

Recent surveys and research studies provide additional insights into the perceptions of AI risks among experts. For instance, a 2022 survey of AI researchers found that the majority believed there is at least a 10% chance that our inability to control AI could cause an existential catastrophe ([Wikipedia](#)). Another survey conducted by AI Impacts in 2023 summarized responses from 2,788 AI researchers regarding the present and future risks of AI, with many expressing concerns about human extinction ([IEEE Spectrum](#)).

## Mitigation Strategies and Future Directions

Given the potential risks associated with AI, various mitigation strategies have been proposed. These include implementing robust safety measures, developing interpretable and steerable AI systems, and promoting transparency and accountability in AI research and development. For example, Anthropic's approach to building reliable AI systems aims to address these concerns and reduce the probability of catastrophic outcomes ([Podtail](#)).

Additionally, some experts advocate for slowing down the pace of AI development to allow for more thorough consideration of the potential risks and ethical implications. Yoshua Bengio, for instance, has decided to push for a slower pace of AI development in an attempt to mitigate the risks ([ABC News](#)).

## Conclusion

While the concept of  $p(\text{doom})$  remains a topic of debate and speculation, it underscores the importance of addressing the potential risks associated with AI. By understanding the factors contributing to  $p(\text{doom})$  and exploring mitigation strategies, researchers and policymakers can work towards ensuring the safe and beneficial development of AI technologies.

# Recent Reports and Expert Opinions on AI Risks

## Expert Opinions on AI Risks

### Probability Estimates and p(doom)

The concept of p(doom), which stands for the probability of AI-induced catastrophe, has gained significant attention among AI researchers and industry leaders. This term encapsulates the likelihood that AI could lead to catastrophic outcomes, including human extinction. For instance, Professor Yoshua Bengio, a prominent figure in AI research, has estimated a 20% probability that AI could result in catastrophic outcomes ([ABC News](#)). This estimate is based on a 50% probability that AI will reach human-level capabilities within a decade and a greater than 50% likelihood that AI or humans will misuse the technology.

Similarly, Geoffrey Hinton, often referred to as the "Godfather of AI," has expressed a 10% chance that AI will lead to human extinction within the next 30 years ([Business Insider](#)). These estimates highlight the varying degrees of concern among experts regarding the potential risks posed by advanced AI systems.

### Recent Surveys and Reports

A recent survey conducted by AI Impacts, involving 2,788 AI researchers, revealed that a significant portion of the community is concerned about the existential risks posed by AI ([IEEE Spectrum](#)). The survey results indicate that many researchers believe there is a non-negligible risk of AI leading to human extinction or other catastrophic outcomes.

The US State Department commissioned a report by AI startup Gladstone, which assessed the risks of AI weaponization and loss of control. The report concluded that advanced AI systems could pose an "extinction-level threat" to humanity ([CNN](#)). The report's findings were based on interviews with over 200 experts, including top executives from leading AI companies, cybersecurity researchers, and national security officials.

# Potential Catastrophic Scenarios

## Biological Warfare and Cyber Attacks

One of the most alarming scenarios discussed in the context of AI risks is the potential for AI to be used in biological warfare. AI systems could be employed to design and deploy bioweapons, leading to mass casualties and global destabilization. The Gladstone report highlighted this risk, noting that AI systems developed in the next 12 to 36 months might be capable of assisting in bioweapon design ([Business Insider](#)).

Cyber attacks are another significant concern. AI systems could execute catastrophic malware attacks, disrupting critical infrastructure and causing widespread chaos. The Gladstone report emphasized the high risk of such attacks, which could lead to mass-casualty events or global destabilization ([CNN](#)).

## Autonomous Weapons and Misinformation

The development of autonomous weapons is another potential risk associated with advanced AI. These weapons could operate without human intervention, making decisions that could lead to unintended and catastrophic consequences. The RAND Corporation's panel discussion on AI risks highlighted the potential for AI to undermine democracy and amplify human stupidity ([The Debrief](#)).

Misinformation and disinformation campaigns are also significant concerns. AI systems could be used to create and spread false information, manipulating public opinion and destabilizing societies. The Gladstone report noted that AI could be used in disinformation campaigns, further exacerbating the risks posed by advanced AI systems ([Business Insider](#)).

# Mitigation Strategies and Policy Recommendations

## Regulation and Oversight

To mitigate the risks associated with advanced AI, experts have proposed various regulatory and oversight measures. Some have suggested regulating the computer chips necessary for advanced AI systems, similar to how fissile uranium is regulated, with an international registry and surprise inspections ([The New Yorker](#)). This approach aims to prevent the uncontrolled proliferation of advanced AI technologies.

The Gladstone report recommended creating regulations to slow down the AI race and establishing an AI safety task force to improve AI capabilities ([CNN](#)). However, some experts, like Artur Kiulian, believe that regulation should encourage innovation rather than stifle it. Kiulian argued that creating a task force would be too expensive and that international safeguards might not be effective, as other countries could continue to advance AI without regard for regulation ([Business Insider](#)).

## **Research and Community Resilience**

Independent, high-quality research is crucial for assessing AI's short- and long-term risks and shaping public policy accordingly. Experts have emphasized the need for rigorous research and policy interventions to ensure that AI's integration into society enhances rather than diminishes human well-being ([The Debrief](#)).

Building resilient communities that can withstand the broad spectrum of crises posed by AI is also essential. This involves not only technical fixes but also a broader view of how AI interacts with human systems, such as criminal justice, education, and employment. Dr. Benjamin Boudreaux, a policy researcher, highlighted the importance of looking at the context in which AI is integrated to understand its impact on human well-being ([The Debrief](#)).

## **Diverging Opinions and Debates**

### **Skepticism and Optimism**

While many experts express significant concerns about the risks posed by advanced AI, others are more skeptical. Yann LeCun, another prominent figure in AI research, has described the notion of AI-induced catastrophe as "preposterously ridiculous" ([ABC News](#)). LeCun and others argue that the positive potential of AI to improve the world outweighs the risks of it becoming hostile.

Similarly, Lorenzo Thione, an AI investor, disagrees with the alarmist logic in some reports, arguing that limiting research and advancement would be ineffective and stifle innovation ([Business Insider](#)). Thione believes that the progression of AI capabilities does not necessarily mean that computational power will continue to increase exponentially.



## **Incremental Harms and Long-Term Risks**

Some experts, like Dr. Benjamin Boudreaux, are more concerned about the incremental harms posed by AI rather than sudden, catastrophic events. Boudreaux likens AI to a slow-moving catastrophe, similar to climate change, that could diminish the institutions and agency needed for meaningful human lives ([The Debrief](#)).

Dr. Roman V. Yampolskiy, an AI safety expert, has warned that there is no evidence that AI superintelligence can be safely controlled. Yampolskiy cautions that without proof of control, AI development should be halted to prevent potential existential catastrophes ([The Debrief](#)).

In summary, the debate over AI risks is complex and multifaceted, with experts offering varying opinions on the likelihood and nature of potential catastrophic outcomes. While some advocate for stringent regulation and oversight, others emphasize the need for innovation and community resilience. The ongoing discourse highlights the importance of continued research and informed policymaking to navigate the challenges posed by advanced AI.

## **Potential Scenarios and Mitigation Strategies**

### **Potential Scenarios of AI-Induced Catastrophes**

#### **Full Automation of Labor and Economic Disruption**

The concept of "full automation of labor" (FAOL) is a scenario where AI systems take over all human jobs, leading to massive economic disruption. According to a survey by AI Impacts, a significant portion of AI researchers believe that FAOL could occur within the next few decades ([source](#)). This scenario raises concerns about widespread unemployment, economic inequality, and social unrest. The transition to a fully automated economy would require substantial changes in social safety nets, education systems, and economic policies to mitigate the negative impacts on society.

#### **Autonomous Weapons and Military Conflicts**

The development of autonomous weapons is another potential scenario that poses a significant risk. AI-driven weapons systems could be used in military conflicts,

leading to unintended escalations and loss of human control over warfare. The [Business Insider report](#) commissioned by the US State Department highlights the risk of AI weaponization, including biowarfare, mass cyber-attacks, and autonomous robots. The report emphasizes the need for international regulations to prevent the proliferation of autonomous weapons and ensure that AI systems are used responsibly in military applications.

## **Cybersecurity Threats and Infrastructure Attacks**

AI systems could be exploited to conduct large-scale cyber-attacks, targeting critical infrastructure such as power grids, financial systems, and communication networks. The [Gladstone report](#) identifies cyber-attacks as one of the highest risks associated with advanced AI. AI-driven malware and hacking tools could bypass traditional security measures, leading to catastrophic consequences. To mitigate this risk, it is essential to develop robust cybersecurity frameworks and invest in AI-driven defense mechanisms.

## **Disinformation and Social Manipulation**

AI systems can be used to create and spread disinformation, manipulate public opinion, and interfere with democratic processes. The [MIT Technology Review](#) article discusses the potential for AI to be used in disinformation campaigns and election interference. AI-generated deepfakes and social media bots can amplify false information, erode trust in institutions, and destabilize societies. Addressing this threat requires collaboration between governments, tech companies, and civil society to develop strategies for detecting and countering disinformation.

## **Mitigation Strategies for AI Risks**

### **Regulation and Oversight**

One of the primary strategies for mitigating AI risks is the implementation of robust regulatory frameworks. The [New Yorker article](#) suggests that AI technologies should be regulated similarly to nuclear materials, with international registries and surprise inspections. This approach would ensure that AI development is closely monitored and that safety measures are in place to prevent misuse. Additionally, the establishment of an international agency, akin to the UN's nuclear watchdog, could oversee AI development and enforce compliance with global standards.

## **Research and Development of Safe AI**

Investing in research focused on AI safety and alignment is crucial for mitigating risks. Organizations like Anthropic have pledged to prioritize safety in their AI development processes ([source](#)). Developing AI systems that align with human values and ethical principles can reduce the likelihood of unintended harmful outcomes. This includes creating AI models that are transparent, interpretable, and capable of being controlled by human operators.

## **Scenario Planning and Risk Assessment**

Scenario planning is a valuable tool for anticipating and preparing for potential AI risks. The [ARL-CNI report](#) outlines various scenarios, ranging from the failure of AI applications to the emergence of superhuman AI capabilities. By considering a wide range of plausible futures, policymakers and stakeholders can develop strategies to address different risk scenarios. This approach helps in identifying critical uncertainties and crafting policies that are flexible and adaptive to changing circumstances.

## **International Collaboration and Governance**

Addressing AI risks requires global cooperation and governance. The [Scientific American article](#) emphasizes the importance of international collaboration in mitigating AI risks. Establishing global agreements on AI development, usage, and safety standards can prevent competitive races that prioritize speed over safety. Collaborative efforts can also facilitate the sharing of best practices, research findings, and resources to enhance AI safety on a global scale.

## **Public Awareness and Education**

Raising public awareness about AI risks and promoting education on AI ethics and safety is essential for building a resilient society. The [TechXplore article](#) highlights the need for policymakers to address the subtle but significant impacts of AI on society. Educating the public about the potential risks and benefits of AI can foster informed decision-making and encourage responsible use of AI technologies. Additionally, promoting interdisciplinary research and dialogue between AI experts, ethicists, and social scientists can contribute to a more comprehensive understanding of AI risks and mitigation strategies.

## Conclusion

The potential scenarios of AI-induced catastrophes and the corresponding mitigation strategies underscore the complexity and urgency of addressing AI risks. By implementing robust regulations, investing in safe AI research, engaging in scenario planning, fostering international collaboration, and raising public awareness, society can navigate the challenges posed by advanced AI technologies. These efforts are crucial for ensuring that AI development aligns with human values and contributes to the well-being and safety of humanity.

## References

- <https://www.scientificamerican.com/article/ai-is-an-existential-threat-just-not-the-way-you-think/?ref=jardel.co>
- <https://rollcall.com/2024/06/11/downplaying-ais-existential-risks-is-a-fatal-error-some-say/>
- <https://spectrum.ieee.org/ai-existential-risk-survey>
- <https://www.fastcompany.com/90994526/pdoom-explained-how-to-calculate-your-score-on-ai-apocalypse-metric>
- <https://www.newyorker.com/magazine/2024/03/18/among-the-ai-doomsayers>
- <https://www.businessinsider.com/ai-report-risks-human-extinction-state-department-expert-reaction-2024-3?op=1>
- <https://www.abc.net.au/news/2023-07-15/whats-your-pdoom-ai-researchers-worry-catastrophe/102591340>
- <https://www.forbes.com/sites/forbesbusinesscouncil/2024/06/05/the-risks-and-rewards-of-ai-strategies-for-mitigation-and-containment/>
- <https://techxplore.com/news/2023-07-ai-existential-threatjust.html>
- <https://www.arl.org/wp-content/uploads/2024/05/ARL-CNI-2035-AI-Scenarios-5May2024.pdf>
- <https://www.forbes.com/sites/forbestechcouncil/2024/02/02/whats-next-for-ai-the-next-wave-of-use-cases-in-2024-and-beyond/>
- <https://www.technologyreview.com/2023/06/06/1074077/to-avoid-ai-doom-learn-from-nuclear-safety/>
- <https://www.technologyreview.com/2024/01/05/1086203/whats-next-ai-regulation-2024/>