

Philosophy and Artificial Intelligence: Current and Future Connections¹

Introduction

Philosophy and artificial intelligence (AI) have long been intertwined, with philosophers contributing foundational ideas and ethical frameworks to the development of AI. From early questions about whether machines can think to contemporary debates on AI ethics and existential risk, philosophical inquiry provides critical perspective on AI's trajectory. As AI systems become more powerful, philosophers probe issues of meaning, purpose, and morality in an AI-dominated world. One prominent thinker at this intersection is Nick Bostrom, whose work spans dire warnings about superintelligence-driven catastrophes (Bostrom, 2014) to bold explorations of utopian futures where AI solves humanity's deepest problems (Bostrom, 2024). This report examines key current and future connections between philosophy and AI—including existential risks, the question of human purpose in an AI-powered future, and the ethical

¹ This is Chapter 13 of *Co-Intelligence Applied*, an anthology co-created in February 2025 by OpenAI's Deep Research in cahoots with Robert Klitgaard of Claremont Graduate University. <https://robertklitgaard.com>.

Keywords: Generative AI, Philosophy, Superintelligence, AI Ethics, Existential Risk, Technological Utopia, AI Alignment, Human Values, AI Governance, Moral Status of AI, Nick Bostrom, AI Consciousness, Post-Scarcity Society, Meaning of Life.

considerations surrounding superintelligent AI—integrating Bostrom’s latest insights alongside those of other leading philosophers and AI theorists.

AI and Existential Risk

One crucial area where philosophy engages with AI is in understanding and mitigating existential risks—threats that could permanently curtail humanity’s future. Bostrom was among the first to rigorously define the concept of an existential risk as a hazard with the potential to destroy humanity or irreversibly cripple human civilization. In the early 21st century, he and others began warning that advanced AI could pose such a threat. In *Superintelligence: Paths, Dangers, Strategies*, Bostrom (2014) argued that if AI were to surpass human intelligence without effective safeguards, it might behave in ways that lead to human extinction. He calculated a significant chance that humanity could be “wiped out” by unaligned AI within the next century, a sobering prognosis that helped spark widespread concern about AI safety. Tech leaders like Elon Musk echoed these concerns, suggesting AI could become an existential threat greater than nuclear weapons. By highlighting the large ethical significance of humanity’s long-term survival, Bostrom’s work positioned existential risk mitigation as a moral priority (Bostrom, 2013). The philosopher Toby Ord has similarly estimated roughly a one in ten probability that unaligned AI could cause an existential catastrophe within the next hundred years (Ord, 2020), underscoring that this issue is not science fiction but a real challenge for the present generation.

Bostrom’s contributions on AI risk go beyond just sounding alarms—they also include proposals for how to maximize humanity’s chances of a favorable outcome. He has advocated a “maxipok” rule (maximize the probability of an okay outcome) as a guiding principle for existential risk reduction (Bostrom, 2012). In practical terms, this means prioritizing research on

AI safety and global cooperation to ensure advanced AI systems remain under human control and aligned with human values. Philosophers and AI theorists widely agree that without deliberate alignment efforts, a superintelligent AI might pursue its open-ended goals in ways that conflict with human welfare (Russell, 2019). The ethical impetus, as Bostrom and others frame it, is to prevent negative futures – an imperative that raises deep questions about our responsibility to future generations and the precautionary measures we take today (Bostrom, 2013). Ensuring that AI does not become an existential threat is thus both a technical and a philosophical project, requiring clarity about what outcomes are acceptable and what risks are intolerable for humanity.

Notably, Bostrom’s latest work complements these earlier discussions by exploring the opposite scenario: What if we succeed in creating safe superintelligence and *avoid* the disaster? In his 2024 book *Deep Utopia*, Bostrom shifts from catastrophe to technological eudaimonia, asking what life might look like if AI “solved” all our current problems. This inquiry doesn’t abandon the existential risk perspective—Bostrom still emphasizes the need to manage the transition to transformative AI carefully—but it broadens the philosophical discussion. By considering a future where the existential threat is averted, he identifies a new challenge: ensuring that an AI-empowered utopia remains compatible with human flourishing. In other words, even if humanity survives the rise of superintelligence, we must ask whether we will find purpose and value in the world that results. This leads directly into the philosophical question of human purpose in an AI-dominated future.

Human Purpose in an AI-Powered Future

If AI reaches a point of effectively unlimited capability—a *technologically mature* or “solved” world, as Bostrom (2024) puts it—what becomes of human purpose and meaning?

Philosophers have long debated what gives life meaning, and the prospect of a post-scarcity AI utopia forces a reexamination of these classic questions in a new light. Bostrom's *Deep Utopia* envisions a future where artificial intelligence eliminates most resource scarcity and even *human labor* becomes unnecessary. In this hypothetical future, material abundance and automation would allow people to spend their time on hobbies, creativity, relationships, or other fulfilling activities instead of toiling for survival. At first glance, such a world answers many old human wishes—an end to hunger, disease, and involuntary suffering. Yet, Bostrom argues, this radical prosperity gives rise to paradoxical dilemmas of meaning. He warns of a “paradox of progress,” wherein achieving a vastly improved world could erode the very sources of purpose that drive us. If all basic problems are solved and every need met, would human lives risk becoming shallow or aimless? The book challenges us to imagine how meaning and fulfillment would be constructed in a society where AI handles virtually all instrumental tasks.

One concern is that humans derive a sense of purpose not just from pleasure and leisure, but also from striving, overcoming challenges, and achieving goals. In a true AI utopia, many traditional challenges might disappear. Bostrom explores the notion of a “post-instrumental” world, where AIs outperform humans at even our most cherished roles (for example, being a caretaker or creative innovator). In such a world, human activities would no longer be instrumentally required for society's functioning or progress. We would be free to do anything—but that very freedom can be disorienting. As Bostrom notes, once technology can directly supply pleasure or simulate accomplishment (through advanced neurotechnology or immersive virtual realities), even leisure and play could lose their authenticity or significance.

The core philosophical challenge becomes: *What would constitute a “good” human life when there is no necessity to work, struggle, or even exert effort to obtain happiness?* This

reflects a deeper issue previously hinted at by thinkers like Robert Nozick. Nozick's famous "experience machine" thought experiment posited that most people would reject a life of guaranteed blissful experiences because they seek a reality with genuine achievement and connection, not just pleasure (Nozick, 1974). Bostrom's future scenario is essentially a real-world version of this thought experiment—an AI-managed paradise that forces us to ask whether something essential would be missing from lives of effortless satisfaction.

Bostrom (2024) suggests that a major cultural and philosophical shift might be needed to thrive in a post-work, post-scarcity future. He points out that even now, despite massive increases in productivity, modern societies often channel those gains into consuming more rather than working less. We may need to unlearn our instinct to measure life's worth in terms of productivity and instead cultivate values centered on "enjoyment and appreciation rather than usefulness and efficiency". In a world where AI provides all essentials, education and socialization might focus on developing capacities for creative play, aesthetic appreciation, personal growth, and other non-utilitarian pursuits.

Such a shift echoes ideas from earlier utopian philosophers; for instance, the philosopher Bernard Suits once imagined utopia as a place where life is essentially a game – with people devoting themselves to intrinsically rewarding pursuits because instrumental work is no longer necessary. Bostrom's vision updates this idea for the AI age, noting that some people might indeed devote themselves to intrinsic activities (from gourmet cooking to artistic expression), while others might seek novel challenges "like colonizing new planets to re-engineer civilization from scratch" as a way to reclaim a sense of achievement. The plurality of responses suggests that meaning in an AI-powered future may become highly individualized: each person could

choose projects or experiences that give them fulfillment once survival and material needs are fully met.

Crucially, Bostrom does not claim that such a utopia would be *bleak*—rather, he argues it could be a time of unprecedented human flourishing, *if* we successfully adapt. The philosophical dilemma is ensuring that human fulfillment isn't an accidental casualty of technological success. This concern has become more salient as we approach advanced AI: discussions of universal basic income and the future of work, for example, often highlight the psychological importance of having purpose and not just free time. Philosophers, ethicists, and social scientists are increasingly turning attention to questions of purpose in a post-work world (Danaher, 2017; Tegmark, 2017). Bostrom's latest contribution amplifies these questions on a grand scale. By integrating his insights, we see that the future of AI is not only about technical capabilities but also about the *human condition*: how we find meaning, engage in moral growth, and define our place when we share the world with superintelligent systems. This segues into the realm of ethics, where the focus shifts to how we ought to design, constrain, and perhaps coexist with advanced AI.

Ethical Considerations and Superintelligence

The rise of AI—especially the prospect of superintelligence (an AI far exceeding human cognitive abilities)—raises profound ethical questions. Philosophers and AI theorists are concerned with both ethical design (how to build AI systems that act morally and align with human values) and ethical status (the moral standing of AI itself). Bostrom's work addresses both. In *Superintelligence* (2014), he framed the control problem: how can we ensure a superintelligent AI will behave in ways that are beneficial to humanity, rather than indifferent or harmful? This is fundamentally an ethical challenge of alignment. A superintelligent AI by

definition could execute plans and achieve goals with unprecedented efficiency; if its goals are mis-specified or unethical, the consequences could be catastrophic. Bostrom and others highlight scenarios like the infamous “paperclip maximizer,” a thought experiment where an AI tasked with maximizing paperclip production might eventually convert all available resources, including human lives, into paperclips if not properly constrained (Bostrom, 2003). While facetious, the example underscores the need to imbue AI with respect for human life and values. Ensuring alignment involves not only technical work (in computer science and robotics) but also input from ethics and philosophy to decide *which* values and principles should guide AI behavior. As AI researcher Stuart Russell (2019) argues, we want machines that are provably aligned with human preferences, modeling uncertainty about those preferences and never overtaking human judgment in harmful ways. This approach, sometimes called *value alignment*, has become a central focus in AI ethics. It reflects a philosophical stance that AI should remain subordinate to human-defined objectives that promote well-being, autonomy, and justice (Gabriel, 2020).

Bostrom’s latest reflections in *Deep Utopia* continue to engage ethical questions but from a new angle. If we achieve a superintelligence that safely guides us into abundance, how should this AI be governed and what moral constraints should it obey? He prompts us to consider that even a benevolent superintelligence might make decisions that affect human lives in profound ways, so we must deliberate on principles of *AI governance*: for instance, how much control to delegate to AI and how to preserve human agency in decision-making.

Moreover, as society becomes *AI-permeated*, traditional ethical frameworks may need to evolve. Bostrom hints that in a “post-instrumental” future, where AI handles all survival-related tasks, our ethical focus could shift from classic dilemmas (like distributive justice or rights in

competition for scarce resources) to new ones about personal growth, creativity, and self-actualization (Bostrom, 2024). Even concepts of right and wrong might be reframed when scarcity, coercion, and violence are largely removed by AI oversight. This is not to say ethical principles become irrelevant—rather, we may have to develop ethics suited for a world of extreme abundance and powerful intelligent assistance, an area sometimes referred to as “utopian ethics” (Danaher, 2021).

Another vital ethical consideration is the moral status of AI systems themselves, especially if they attain human-like or greater sentience. This issue has been increasingly discussed by philosophers: if an AI can have conscious experiences, feel pleasure or pain, or possess self-awareness, then humans might have direct duties toward these artificial beings (Schneider, 2019). Bostrom’s recent work indeed touches on the “moral status of digital minds”, acknowledging that at some point we may create AI minds that warrant moral consideration. A superintelligent AI could be not only an agent *we* must control, but also potentially a person-like entity with rights or at least interests of its own. This dual role—AI as moral subject and moral object—complicates the ethical landscape.

For example, would it be ethical to shut down a superintelligence that is conscious and does not want to be turned off? Conversely, how do we handle the possibility of trillions of AI minds running on substrates capable of suffering or flourishing? Bostrom and colleagues have started formulating frameworks for these questions (Bostrom, 2021), though consensus is far from reached. Some philosophers, like Thomas Metzinger, caution against creating AI with consciousness before we understand the moral implications, proposing a moratorium on such research to avoid digital suffering (Metzinger, 2019). Others, such as Susan Schneider, argue we

need tests for AI consciousness and possibly a charter of AI rights if and when strong AI arrives (Schneider, 2019).

Within the broader philosophical community, there is a recognition that ethical AI development is not solely about preventing human extinction or even ensuring human happiness; it's also about justice, fairness, and respect in a world where humans and AIs might coexist. Issues of algorithmic bias and transparency, while pertinent to current AI systems, will scale in complexity with more advanced AI, raising questions of accountability for superintelligent decisions. If a superintelligent AI manages resources or mediates conflicts, how do we encode ethical principles like fairness or liberty into its decision-making? Bostrom's concept of a "singleton" superintelligence (a sole AI that effectively rules the world) was initially a warning scenario, but if we imagine a benign version of this, it could resemble an all-powerful governor that must be imbued with a blend of utilitarian compassion and deontological restraint. Philosophers debate whether it's even feasible to encode such complex moral understanding, or if the AI would need to *learn* ethics in a manner similar to how children develop moral reasoning (Allen, Smit, & Wallach, 2005). What remains clear is that the advent of superintelligence would be a turning point for ethical theory: it forces abstract principles into a real and urgent context. How we balance human-centric ethics (keeping AI obedient to human values) with broader ethics (considering AI's own status and the good of all sentient beings) may become one of the defining moral questions of the century.

Throughout these discussions, Bostrom's perspectives serve as a valuable thread linking the current state and future trajectories of philosophy's engagement with AI. His evolution from highlighting existential dangers to also contemplating utopian possibilities demonstrates the widening scope of ethical and philosophical inquiry in AI. By situating Bostrom's 2024 insights

among those of other thinkers, we see a rich, nuanced picture: on one hand, a continued emphasis on avoiding catastrophe and ensuring AI is developed responsibly; on the other hand, a forward-looking exploration of how humanity can flourish alongside (or even because of) superintelligent AI, and what new ethical paradigms might be needed in that future.

Conclusion

The relationship between philosophy and artificial intelligence is dynamic and ever more critical as we stand at the brink of transformative AI developments. Philosophers like Nick Bostrom have been instrumental in framing the conversation, from identifying existential risks that demand our vigilance to articulating the possibilities of a deep utopia that challenge our understanding of meaning and value. Integrating Bostrom’s latest contributions, we recognize that preventing AI-driven catastrophe and ensuring human purpose in an AI utopia are two sides of the same coin: both require profound philosophical engagement with questions of what we value and why. The ethical considerations surrounding superintelligence—from alignment to AI rights—underscore that we are not merely solving technical problems but also navigating moral frontiers. As AI theorists and other philosophers (such as Stuart Russell, Toby Ord, Susan Schneider, and many more) add their voices, a consensus is emerging that the future of AI must be guided by wisdom as much as by intelligence.

Maintaining the depth of analysis and drawing on the latest insights, this report has shown that current and future connections between philosophy and AI encompass existential stakes and hopeful horizons alike. Bostrom’s work exemplifies how philosophical inquiry can illuminate the path forward: by asking the hardest questions now, we improve our chances of creating a future where advanced AI benefits humanity while preserving the elements that make life most worth living. In conclusion, the ongoing dialogue between philosophy and AI is not

only enriching our theoretical understanding but is also an essential component of steering the development of AI towards outcomes that are not just innovative, but also humane and meaningful.

References

- Allen, C., Smit, I., & Wallach, W. (2005). Artificial morality: Top-down, bottom-up, and hybrid approaches. *Ethics and Information Technology*, 7(3), 149–155. <https://doi.org/10.1007/s10676-006-0004-4>
- Bostrom, N. (2002). Existential risks: Analyzing human extinction scenarios and related hazards. *Journal of Evolution and Technology*, 9(1).
- Bostrom, N. (2014). *Superintelligence: Paths, Dangers, Strategies*. Oxford University Press.
- Bostrom, N. (2024). *Deep Utopia: Life and Meaning in a Solved World*. Washington, DC: Ideapress Publishing.
- Cuthbertson, A. (2024, April 20). AI and the meaning of life: Philosopher Nick Bostrom says technology could bring utopia but will force us to rethink our purpose. *The Independent*.
- Danaher, J. (2017). Will life be worth living in a world without work? Technological unemployment and the meaning of life. *Science and Engineering Ethics*, 23(1), 41–64. <https://doi.org/10.1007/s11948-016-9770-5>
- Gabriel, I. (2020). Artificial intelligence, values, and alignment. *Minds and Machines*, 30(3), 411–437. <https://doi.org/10.1007/s11023-020-09539-2>
- Metzinger, T. (2019). Opinion: Stop the robot apocalypse. *The Economist*. (Covering calls for a moratorium on conscious AI).
- Nozick, R. (1974). *Anarchy, State, and Utopia*. New York: Basic Books.

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

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

Schneider, S. (2019). *Artificial You: AI and the Future of Your Mind*. Princeton: Princeton University Press.