

The intelligence in most of its dimensions,

[Business](#), [Decision Making](#)



The implications of technological innovation for sustainability are becoming increasingly complex with information technology moving machines from being mere tools for production or objects of consumption to playing a role in economic decision making. This emerging role will acquire overwhelming importance if artificial intelligence is underway to outperform human intelligence in most of its dimensions, thus becoming superintelligence. Superintelligence explores the future of artificial intelligence and related technologies and the risks they may pose to human civilization. The author addresses the topic of superintelligence in the field of artificial intelligence in both a very captivating structure and an interesting approach. Absent the optimal background, the most important arguments of the book, will still be accessible, though some of the technical vocabulary may be a source of frustration for some. The book skillfully demonstrates the potential for serious thinking aimed at the long-term future.

A very interesting thought that the author discussed is superintelligence is compatible with almost any final goal. This thought is also scary somehow. This was introduced in the book as the orthogonality thesis and has a crucial consequence that the possession of superintelligence does not imply being wise or benevolent. In the matter of fact, there is no reason to believe that a superintelligent being has high moral standards; Even a non-malevolent superintelligence may destroy everything dear to us or perform otherwise morally terrible actions.

So if we want the self-preservation of humankind, is it enough to be careful and cautious with artificial intelligence or won't that be enough? Is it already too late and did we already reach the point of no return? Simulating what

<https://assignbuster.com/the-intelligence-in-most-of-its-dimensions/>

human people wish requires simulating human people, according to the book, terminating a simulation could then easily turn out as a genocide. For example, charged with the goal of computing the numeric value of pi, or manufacturing many paperclips, a superintelligent being might proceed by an unbridled acquisition of physical resources, in order to facilitate its computations or manufacturing capacity, thereby appropriating our bodies as a convenient source of atoms, or modifying our environment in a way that results in our extinction. This leads to the question: what if reaching our artificial intelligence goals, whatever these may be, for comfort reasons or science reasons, will our ancestors' survival goal disappear. Another theory mentioned in the book assumes that superintelligences tend to converge to the same goals. One of their reasons is that they may be because most humans have similar values, so they will eventually adapt these values to satisfy humans. This concept is known as instrumental convergence. Although superintelligent beings are offered a wide variety of final goals, they would most of the time pursue the same intermediate goal, whose satisfaction tends to enable the satisfaction of almost any final goal. Is assuming that there is a rationally correct morality just us being optimistic about the future of artificial intelligence or is it realistic that a sufficiently rational artificial intelligence will acquire this proper morality and begin to act according to it? As much as we all want or choose to believe that humans have similar virtues, believe in some common good deeds and work for similar goals, just by looking at how messed up our world was and is right now (at least the way I see it based on the wars humankind witnesses and

the environmental destruction of the planet), it is hard to assume that there always exists an intermediate goal.

A goal that is good for one person, could be disastrous for another. So how can we assume that we can make an efficient real-world artificial intelligence algorithm that is suited for the sake of everyone. Even if a human society were highly motivated to design an efficient real-world algorithm, and were given enough time to do so along with huge amounts of resources, training and knowledge about artificial intelligence, I doubt that it will come a point where everyone could agree upon the same solution for a certain problem. How can we then assume that artificial intelligence systems- if faced with a real problem- will take the right decision that will cause no harm to humans, if humans themselves fail to do this sometimes? So do I believe that superintelligence is realistic? Well, half a century after the first electric computer, we still have nothing that even resembles an intelligent machine, if by 'intelligent' we mean possessing the kind of general-purpose smartness that we humans pride ourselves on. However, neither the fact that machine intelligence would be challenging nor the fact that some past predictions were wrong is a good ground for concluding that artificial intelligence will never be created.

Indeed, to assume that artificial intelligence will take thousands of years to develop seems at least as unwarranted as to make the opposite assumption. Already on the first few pages I sensed the speculative narrative and the intention of the author to hold forth to conclusions such as mankind's

doomsday. Unfortunately, there were many speculations that I found difficult to recognize without turning my basic logic circuits off.

Despite the book's clear explanation of why superintelligent AI may have arbitrarily negative consequences and why it is important to begin addressing the issue well in advance, the author does not base his case on predictions that superhuman AI systems are imminent. He writes, "It is no part of the argument in this book that we are on the threshold of a big breakthrough in artificial intelligence, or that we can predict with any precision when such a development might occur." Pointing to long-term risks from AI is however not equivalent to claiming that superintelligence and its accompanying risks are imminent. From my point of view, I do think that AI presents an existential risk and it should be seriously recognized. This risk does not necessarily rise from spontaneous malicious consciousness, but rather from the unpredictability and potential irreversibility of deploying an optimization process more intelligent than the humans who specified its objectives. However, no one can say for sure if it is 100% happening and when will it happen. We all base our opinions on assumptions depending on how optimistic or not we are and according to our own knowledge of the evolution of technological inventions. But the future is always unpredictable.

One thing is for sure the challenge of superintelligence cannot be separated from the other major environmental and social challenges, demanding a fundamental transformation along the lines of degrowth. With machines outperforming humans in their functions, maybe our social skills will still keep humans ahead of machines. As a conclusion, my general

impression was that this book might not be the most pleasurable read at some sections, sometimes repetitive and not straightforward, nevertheless, sticking with it till the end was for me worth its while. The book has a good enough logical structure, and most chapters end with a summary.

The author succeeded in giving a detailed analysis about this very important theme of the modern world and engaged with the topic very competently. To my viewpoint Bostrom succeeds in arguing that the development of superintelligent machines will, if not properly managed, create catastrophic risks to humanity.