

## ARE WE LIVING IN A COMPUTER SIMULATION?

BY NICK BOSTROM

*I argue that at least one of the following propositions is true: (1) the human species is very likely to become extinct before reaching a 'posthuman' stage; (2) any posthuman civilization is extremely unlikely to run a significant number of simulations of its evolutionary history (or variations thereof); (3) we are almost certainly living in a computer simulation. It follows that the belief that there is a significant chance that we shall one day become posthumans who run ancestor-simulations is false, unless we are currently living in a simulation. I discuss some consequences of this result.*

### I. INTRODUCTION

Many works of science fiction as well as some forecasts by serious technologists and futurologists predict that enormous amounts of computing power will be available in the future. Let us suppose for a moment that these predictions are correct. One thing that later generations might do with their super-powerful computers is run detailed simulations of their forebears or of people like their forebears. Because their computers would be so powerful, they could run a great many such simulations. Suppose that these simulated people are conscious (as they would be if the simulations were sufficiently fine-grained and if a certain quite widely accepted position in the philosophy of mind is correct). Then it could be the case that the vast majority of minds like ours do not belong to the original race but rather to people simulated by the advanced descendants of an original race. It is then possible to argue that if this were the case, we would be rational to think that we are likely to be among the simulated minds rather than among the original biological ones. Therefore if we do not think that we are currently living in a computer simulation, we are not entitled to believe that we shall have descendants who will run lots of simulations of their forebears. That is the basic idea. The rest of this paper will spell it out more carefully.

Apart from the interest this thesis may hold for those engaged in futuristic speculation, there are also more purely theoretical rewards. The argument is a stimulus for formulating some methodological and metaphysical questions,

and it suggests naturalistic analogies of certain traditional religious conceptions, which some may find amusing or thought-provoking.

The structure of the paper is as follows. First, I formulate an assumption which I need to import from the philosophy of mind in order to get the argument started. Secondly, I consider some empirical reasons for thinking that running vastly many simulations of human minds would be within the capability of a future civilization that has developed many of those technologies that can already be shown to be compatible with known physical laws and engineering constraints. This part is not philosophically necessary, but it provides an incentive for paying attention to the rest. Then follows the core of the argument, which makes use of some simple probability theory, and a section providing support for a weak indifference principle the argument employs. Lastly, I discuss some interpretations of the disjunction mentioned in the abstract, which forms the conclusion of the simulation argument.

## II. THE ASSUMPTION OF SUBSTRATE-INDEPENDENCE

A common assumption in the philosophy of mind is that of *substrate-independence*. The idea is that mental states can supervene on any of a broad class of physical substrates. Provided a system implements the right sort of computational structures and processes, it can be associated with conscious experiences. It is not an essential property of consciousness that it is implemented on carbon-based biological neural networks inside a cranium: silicon-based processors in a computer could in principle do the trick too.

Arguments for this thesis have been given in the literature, and although it is not entirely uncontroversial, I shall here take it as given.

The argument I shall present does not, however, depend on any very strong version of functionalism or computationalism. For example, I need not assume that the thesis of substrate-independence is *necessarily* true (either analytically or metaphysically) – merely that a computer running a suitable program would in fact be conscious. Moreover, I need not assume that in order to create a mind on a computer it would be necessary to program it in such a way that it behaves like a human in all situations, including passing the Turing test, etc. I need only the weaker assumption that it would suffice for the generation of subjective experiences that the computational processes of a human brain are structurally replicated in suitably fine-grained detail, such as on the level of individual synapses. This attenuated version of substrate-independence is quite widely accepted.

Neurotransmitters, nerve growth factors and other chemicals that are smaller than a synapse clearly play a role in human cognition and learning.

The substrate-independence thesis is not that the effects of these chemicals are small or irrelevant, but rather that they affect subjective experience only via their direct or indirect influence on computational activities. For example, if there can be no difference in subjective experience without there also being a difference in synaptic discharges, then the requisite detail of simulation is at the synaptic level (or higher).

### III. THE TECHNOLOGICAL LIMITS OF COMPUTATION

At our current stage of technological development, we have neither sufficiently powerful hardware nor the requisite software to create conscious minds in computers. But persuasive arguments have been given to the effect that if technological progress continues unabated, then these technological shortcomings will eventually be overcome. Some authors argue that this stage may be only a few decades away.<sup>1</sup> Yet present purposes require no assumptions about the time-scale. The simulation argument works equally well for those who think that it will take hundreds of thousands of years to reach a 'posthuman' stage of civilization, where humankind has acquired most of the technological capabilities that one can currently show to be consistent with physical laws and with material and energy constraints.

Such a mature stage of technological development will make it possible to convert planets and other astronomical resources into enormously powerful computers. It is currently hard to be confident in any upper bound on the computing power that may be available to posthuman civilizations. As we are still lacking a 'theory of everything', we cannot rule out the possibility that novel physical phenomena, not allowed for in current physical theories, may be utilized to transcend those constraints that in our current understanding impose theoretical limits on the information processing attainable in a given lump of matter.<sup>2</sup> We can with much greater confidence establish *lower* bounds on posthuman computation, by assuming only mechanisms that are already understood. For example, Eric Drexler has outlined a

<sup>1</sup> See, e.g., K.E. Drexler, *Engines of Creation: the Coming Era of Nanotechnology* (London: Fourth Estate, 1985); N. Bostrom, 'How Long Before Superintelligence?', *International Journal of Futures Studies*, 2 (1998); R. Kurzweil, *The Age of Spiritual Machines: When Computers Exceed Human Intelligence* (New York: Viking, 1999); H. Moravec, *Robot: Mere Machine to Transcendent Mind* (Oxford UP, 1999).

<sup>2</sup> I.e., constraints such as the Bremermann–Bekenstein bound and the black hole limit: H.J. Bremermann, 'Minimum Energy Requirements of Information Transfer and Computing', *International Journal of Theoretical Physics*, 21 (1982), pp. 203–17; J.D. Bekenstein, 'Entropy Content and Information Flow in Systems with Limited Energy', *Physical Review*, D 30 (1984), pp. 1669–79; A. Sandberg, 'The Physics of Information Processing Superobjects: the Daily Life among the Jupiter Brains', *Journal of Evolution and Technology*, 5 (1999).

design for a system the size of a sugar cube (excluding cooling and power supply) that would perform  $10^{21}$  instructions per second.<sup>3</sup> Another author gives a rough estimate of  $10^{42}$  operations per second for a computer with a mass of the order of a large planet.<sup>4</sup> (If we could create quantum computers, or learn to build computers out of nuclear matter or plasma, we could push closer to the theoretical limits. Seth Lloyd calculates an upper bound for a 1 kg computer of  $5 \times 10^{50}$  logical operations per second carried out on  $\sim 10^{31}$  bits.<sup>5</sup> However, it suffices for my purposes to use the more conservative estimate that presupposes only currently known design principles.)

The amount of computing power needed to emulate a human mind can likewise be roughly estimated. One estimate, based on how computationally expensive it is to replicate the functionality of a piece of nervous tissue which we have already understood and whose functionality has been replicated *in silice*, namely, contrast enhancement in the retina, yields a figure of  $\sim 10^{14}$  operations per second for the entire human brain.<sup>6</sup> An alternative estimate, based on the number of synapses in the brain and their firing frequency, gives a figure of  $\sim 10^{16}$ – $10^{17}$  operations per second.<sup>7</sup> Conceivably, even more could be required if we want to simulate in detail the internal workings of synapses and dendritic trees. However, it is likely that the human central nervous system has a high degree of redundancy on the microscale to compensate for the unreliability and noisiness of its neuronal components. One would therefore expect a substantial efficiency gain when using more reliable and versatile non-biological processors.

Memory seems to be no more stringent a constraint than processing power.<sup>8</sup> Moreover, since the maximum human sensory bandwidth is  $\sim 10^8$  bits per second, simulating all sensory events incurs a negligible cost compared to simulating the cortical activity. We can therefore use the processing power required to simulate the central nervous system as an estimate of the total computational cost of simulating a human mind.

If the environment is included in the simulation, this will require additional computing power – how much, depends on the scope and granularity of the simulation. Simulating the entire universe down to the quantum level is obviously infeasible, unless radically new physics is discovered. But in order to get a realistic simulation of human experience, much less is needed – only whatever is required to ensure that the simulated humans, interacting

<sup>3</sup> K.E. Drexler, *Nanosystems* (New York: John Wiley & Sons, 1992).

<sup>4</sup> R.J. Bradbury, 'Matrioshka Brains', *working manuscript* (2002), <http://www.aciveos.com/~bradbury/MatrioshkaBrains/MatrioshkaBrains.html>.

<sup>5</sup> S. Lloyd, 'Ultimate Physical Limits to Computation', *Nature*, 406 (31 August 2000), pp. 1047–54.

<sup>6</sup> H. Moravec, *Mind Children* (Harvard UP, 1989).

<sup>7</sup> See my 'How Long before Superintelligence?'

<sup>8</sup> See references in foregoing footnotes.

in normal human ways with their simulated environment, do not notice any irregularities. The microscopic structure of the inside of the Earth can be safely omitted. Distant astronomical objects can have highly compressed representations: verisimilitude need only extend to the narrow band of properties that we can observe from our planet or solar system spacecraft. On the surface of Earth, macroscopic objects in inhabited areas may need to be continuously simulated, but microscopic phenomena could probably be filled in *ad hoc*. What you see through an electron microscope needs to look unsuspecting, but you usually have no way of confirming its coherence with unobserved parts of the microscopic world. Exceptions arise when we deliberately design systems to harness unobserved microscopic phenomena that operate in accordance with known principles to get results we are able to verify independently. The paradigm case of this is a computer. The simulation may therefore need to include continuous representation of computers down to the level of individual logic elements. This presents no problem, since our current computing power is negligible by posthuman standards.

Moreover, a posthuman simulator would have enough computing power to keep track of the detailed belief-states in all human brains at all times. Therefore, when it saw that a human was about to make an observation of the microscopic world, it could fill in sufficient detail in the simulation in the appropriate domain as and where needed. Should any error occur, the director could edit the states of any brains that have become aware of an anomaly before this spoils the simulation. Alternatively, the director could skip back a few seconds and rerun the simulation so as to avoid the problem.

It thus seems plausible that the main computational cost in creating simulations that are indistinguishable from physical reality for human minds in the simulation resides in simulating organic brains down to the neuronal or sub-neuronal level. As we build more and faster computers, the cost of simulating our machines might eventually come to dominate the cost of simulating nervous systems. While it is not possible to get a very exact estimate of the cost of a realistic simulation of human history, we can use  $\sim 10^{33}$ – $10^{36}$  operations as a rough estimate.<sup>9</sup> As we gain more experience with virtual reality, we shall get a better grasp of the computational requirements for making such worlds appear realistic to their visitors. But in any case, even if the estimate is inaccurate by several orders of magnitude, this does not matter much for my argument. I noted that a rough approximation of the computational power of a planetary-mass computer is  $10^{42}$  operations per second, and that assumes only already known nanotechnological designs, which are probably far from optimal. A single such computer could simulate

<sup>9</sup> 100 billion humans  $\times$  50 years/human  $\times$  30 million secs/year  $\times$   $[10^{14}, 10^{17}]$  operations in each human brain per second  $\approx [10^{33}, 10^{36}]$  operations.

the entire mental history of humankind (I shall call this an *ancestor-simulation*) by using less than one millionth of its processing power for one second. A posthuman civilization may eventually build an astronomical number of such computers. I can conclude that the computing power available to a posthuman civilization is sufficient to run a huge number of ancestor-simulations even if it allocates only a very minute fraction of its resources to that purpose. I can draw this conclusion even while leaving a huge substantial margin of error in all our estimates:

Post-human civilizations would have enough computing power to run hugely many ancestor-simulations even while using only a tiny fraction of their resources for that purpose.

#### IV. THE CORE OF THE SIMULATION ARGUMENT

The basic idea of this paper can be expressed roughly as follows: if there were a substantial chance that our civilization will get to the posthuman stage and run many ancestor-simulations, then how come we are not living in such a simulation?

I shall develop this idea into a rigorous argument. I need to introduce the following notation:

- $f_p$ : Fraction of all human-level technological civilizations that survive to reach a posthuman stage
- $\bar{N}$ : Average number of ancestor-simulations run by a posthuman civilization
- $\bar{H}$ : Average number of individuals that have lived in a civilization before it reaches a posthuman stage.

The actual fraction of all observers with human-type experiences that live in simulations is then

$$f_{sim} = \frac{f_p \bar{N} \bar{H}}{(f_p \bar{N} \bar{H}) + \bar{H}}$$

Writing  $f_I$  for the fraction of posthuman civilizations that are interested in running ancestor-simulations (or that contain at least some individuals who are interested in them and have sufficient resources to run a significant number of such simulations), and  $\bar{N}_I$  for the average number of ancestor-simulations run by such interested civilizations, we have

$$\bar{N} = f_I \bar{N}_I$$

and thus

$$F. f_{sim} = \frac{f_p f_I \bar{N}_I}{(f_p f_I \bar{N}_I) + 1}$$

Because of the immense computing power of posthuman civilizations,  $\bar{N}_I$  is extremely large, as I pointed out in the previous section. What (F) shows is that *at least one* of the following three propositions must be true:

1.  $f_p \approx 0$
2.  $f_I \approx 0$
3.  $f_{sim} \approx 1$

## V. A BLAND INDIFFERENCE PRINCIPLE

I can take a further step and conclude that given the truth of (3), one's credence in the hypothesis that one is in a simulation should be close to unity. More generally, if we knew that a fraction  $x$  of all observers with human-type experiences live in simulations, and we have no information to indicate that our own particular experiences are any more or less likely than other human-type experiences to have been implemented *in vivo* rather than *in machina*, then our credence that we are in a simulation should equal  $x$ :

$$S. Cr(\text{SIM} \mid f_{sim} = x) = x.$$

This step is sanctioned by a very weak indifference principle. Two cases need to be distinguished. The first case, which is the easiest, is where all the minds in question are like our own in the sense that they are exactly qualitatively identical with ours: they have exactly the same information and the same experiences as we have. The second case is where the minds are 'like' each other only in the loose sense of being the sort of minds that are typical of human creatures, but where they are qualitatively distinct from one another and each has a distinct set of experiences. I maintain that even in the latter case, where the minds are qualitatively different, the simulation argument still works, provided that we have no information bearing on the question of which of the various minds are simulated and which are implemented biologically.

A detailed defence of a stronger principle, which implies the above stance for both cases as trivial special instances, has been given in the literature.<sup>10</sup> Space does not permit a recapitulation of that defence here, but I can bring

<sup>10</sup> In, e.g., N. Bostrom, 'The Doomsday Argument, Adam and Eve, UN++, and Quantum Joe', *Synthese*, 127 (2001), pp. 359–87; and most fully in my book *Anthropic Bias: Observation Selection Effects in Science and Philosophy* (New York: Routledge, 2002).

out one of the underlying ideas by rehearsing an analogous situation of a more familiar kind. Suppose that  $x\%$  of the population has a certain genetic sequence  $S$  within the part of their DNA commonly designated as 'junk DNA'. Suppose further that there are no manifestations of  $S$  (short of what would turn up in a gene assay) and that there are no known correlations between having  $S$  and any observable characteristic. Then quite clearly, unless one has had one's DNA sequenced, it is rational to assign a credence of  $x\%$  to the hypothesis that one has  $S$ . And this is so quite irrespective of the fact that the people who have  $S$  have qualitatively different minds and experiences from the people who do not have  $S$ . (They are different simply because all humans have different experiences from one another, not because of any known link between  $S$  and what kind of experiences one has.)

The same reasoning holds if  $S$  is not the property of having a certain genetic sequence but instead the property of being in a simulation, assuming only that we have no information that enables us to predict any differences between the experiences of simulated minds and those of the original biological minds.

It should be stressed that the bland indifference principle expressed by (S) prescribes indifference only between hypotheses about which observer one is, when one has no information about which of these observers one is. It does not in general prescribe indifference between hypotheses when one lacks specific information about which of the hypotheses is true. In contrast with Laplacean and other more ambitious principles of indifference, it is therefore immune to Bertrand's paradox and similar predicaments that tend to plague indifference principles of unrestricted scope.

Readers familiar with the doomsday argument<sup>11</sup> may worry that the bland principle of indifference invoked here is the same assumption as is responsible for getting the doomsday argument off the ground, and that the counter-intuitive nature of some of the implications of the latter incriminates or casts doubt on the validity of the former. This is not so. The doomsday argument rests on a *much* stronger and more controversial premise, namely, that one should reason as if one were a random sample from the set of all people who will ever have lived (past, present, and future) *even though we know that we are living in the early twenty-first century* rather than at some point in the distant past or the future. The bland indifference principle, by contrast, applies only to cases where we have no information about which group of people we belong to.

If betting odds provide some guidance to rational belief, it may also be worth pondering that if everybody were to place a bet on whether they are

<sup>11</sup> See, e.g., J. Leslie, 'Is the End of the World Nigh?', *The Philosophical Quarterly*, 40 (1990), pp. 65–72.



in a simulation or not, then if people use the bland principle of indifference, and consequently place their money on being in a simulation if they know that that is where almost all people are, then almost everyone will win their bets. If they bet on *not* being in a simulation, then almost everyone will lose. It seems better that the bland indifference principle should be heeded.

Further, one can consider a sequence of possible situations in which an increasing fraction of all people live in simulations: 98%, 99%, 99.9%, 99.999%, and so on. As one approaches the limiting case in which *everybody* is in a simulation (from which one can *deductively* infer that one is in a simulation oneself), it is plausible to require that the credence one assigns to being in a simulation should gradually approach the limiting case of complete certainty in a matching manner.

## VI. INTERPRETATION

The possibility represented by proposition (1) is fairly straightforward. If (1) is true, then humankind will almost certainly fail to reach a posthuman level; for virtually no species at our level of development become posthuman, and it is hard to see any justification for thinking that our own species will be especially privileged or protected from future disasters. Conditionally on (1), therefore, we must give a high credence to DOOM, the hypothesis that humankind will go extinct before reaching a posthuman level:

$$Cr(\text{DOOM} \mid f_p \approx 1) \approx 1$$

One can imagine hypothetical situations where we have such evidence as would trump knowledge of  $f_p$ . For example, if we discovered that we were about to be hit by a giant asteroid, this might suggest that we had been exceptionally unlucky. We could then assign a credence to DOOM larger than our expectation of the fraction of human-level civilizations that fail to reach posthumanity. In the actual case, however, we seem to lack evidence for thinking that we are special in this regard, for better or worse.

Proposition (1) does not by itself imply that we are likely to go extinct soon, only that we are unlikely to reach a posthuman stage. This possibility is compatible with our remaining at, or somewhat above, our current level of technological development for a long time before going extinct. Another way for (1) to be true is if it is likely that technological civilization will collapse. Primitive human societies might then remain on Earth indefinitely.

There are many ways in which humanity could become extinct before reaching posthumanity. Perhaps the most natural interpretation of (1) is that we are likely to go extinct as a result of the development of some powerful

but dangerous technology.<sup>12</sup> One candidate is molecular nanotechnology, which in its mature stage would enable the construction of self-replicating nanobots capable of feeding on dirt and organic matter – a kind of mechanical bacteria. Such nanobots, designed for malicious ends, could cause the extinction of all life on our planet.<sup>13</sup>

The second alternative in the simulation argument's conclusion is that the fraction of posthuman civilizations that are interested in running ancestor-simulations is negligibly small. In order for (2) to be true, there must be a strong *convergence* among the courses of advanced civilizations. If the number of ancestor-simulations created by the interested civilizations is extremely large, the rarity of such civilizations must be correspondingly extreme. Virtually no posthuman civilizations decide to use their resources to run large numbers of ancestor-simulations. Furthermore, virtually all posthuman civilizations lack individuals who have sufficient resources and interest to run ancestor-simulations; or else they have reliably enforced laws that prevent such individuals from acting on their desires.

What force could bring about such a convergence? One might speculate that advanced civilizations all develop along a trajectory that leads to recognition of an ethical prohibition against running ancestor-simulations because of the suffering that is inflicted on the inhabitants of the simulation. However, from our present point of view, it is not clear that creating a human race is immoral. On the contrary, we tend to view the existence of our race as constituting a great ethical value. Moreover, convergence on an ethical view of the immorality of running ancestor-simulations is not enough: it must be combined with convergence on a civilization-wide social structure that enables activities considered immoral to be effectively banned.

Another possible convergence point is that almost all individual posthumans in virtually all posthuman civilizations develop in a direction where they lose their desire to run ancestor-simulations. This would require significant changes to the motives driving their human predecessors, for there are certainly many humans who would like to run ancestor-simulations if they could afford to do so. But perhaps many of our human desires will be regarded as silly by anyone who becomes a posthuman. Maybe the scientific value of ancestor-simulations to a posthuman civilization is negligible (which is not too implausible given its unfathomable intellectual superiority), and maybe posthumans regard recreational activities as merely

<sup>12</sup> See my 'Existential Risks: Analyzing Human Extinction Scenarios and Related Hazards', *Journal of Evolution and Technology*, 9 (2001), for a survey and analysis of the present and anticipated future threats to human survival.

<sup>13</sup> See, e.g., Drexler; and R.A. Freitas Jr, 'Some Limits to Global Ecophagy by Biovorous Nanoreplicators, with Public Policy Recommendations', *Zyvex preprint*, April (2000), <http://www.foresight.org/NanoRev/Ecophagy.html>.

a very inefficient way of getting pleasure – which can be obtained much more cheaply by direct stimulation of the brain's reward centres. One conclusion that follows from (2) is that posthuman societies will be very different from human societies: they will not contain relatively wealthy independent agents who have the full gamut of human-like desires and are free to act on them.

The possibility expressed by alternative (3) is conceptually the most intriguing one. If we are living in a simulation, then the cosmos we are observing is just a tiny piece of the totality of physical existence. The physics in the universe where the computer running the simulation is situated may or may not resemble the physics of the world we observe. While the world we see is in some sense 'real', it is not located at the fundamental level of reality.

It may be possible for simulated civilizations to become posthuman. They may then run their own ancestor-simulations on powerful computers they build in their simulated universe. Such computers would be 'virtual machines', a familiar concept in computer science. (Javascript web-applets, for instance, run on a virtual machine – a simulated computer – inside a desktop.) Virtual machines can be stacked: it is possible to simulate one machine simulating another machine, and so on, in arbitrarily many steps of iteration. If we do go on to create our own ancestor-simulations, then since this would be strong evidence against (1) and (2), we would therefore have to conclude that we live in a simulation. Moreover, we would have to suspect that the posthumans running our simulation are themselves simulated beings; and their creators in turn may also be simulated beings.

Reality may thus contain many levels. Even if it is necessary for the hierarchy to bottom out at some stage – the metaphysical status of this claim is somewhat obscure – there may be room for a large number of levels of reality, and the number could be increasing over time. (One consideration that counts against the multi-level hypothesis is that the computational cost for the basement-level simulators would be very great. Simulating even a single posthuman civilization might be prohibitively expensive. If so, then we should expect our simulation to be terminated when we are about to become posthuman.)

Although all the elements of such a system can be naturalistic, even physical, it is possible to draw some loose analogies with religious conceptions of the world. In some ways, the posthumans running a simulation are like gods in relation to the people inhabiting the simulation: the posthumans created the world we see; they are of superior intelligence; they are 'omnipotent' in the sense that they can interfere in the workings of our world even in ways that violate its physical laws; and they are 'omniscient' in the sense that they can monitor everything that happens. However, all the

demigods except those at the fundamental level of reality are subject to sanctions by the more powerful gods living at lower levels.

Further rumination on these themes could climax in a *naturalistic theogony* that would study the structure of this hierarchy, and the constraints imposed on its inhabitants by the possibility that their actions on their own level may affect the treatment they receive from dwellers of deeper levels. For example, if nobody can be sure that they are at the basement-level, then everybody would have to consider the possibility that their actions will be rewarded or punished, perhaps using moral criteria, by their simulators. An afterlife would be a real possibility. Because of this fundamental uncertainty, even the basement civilization may have a reason to behave ethically. The fact that it has such a reason for moral behaviour would of course add to everybody else's reason for behaving morally, and so on, in a truly virtuous circle. One might get a kind of universal ethical imperative, which it would be in everybody's self-interest to obey, as it were, 'from nowhere'.

In addition to ancestor-simulations, one may also consider the possibility of more selective simulations that include only a small group of humans or a single individual. The rest of humanity would then be zombies or 'shadow-people' – humans simulated only at a level sufficient for the fully simulated ones not to notice anything suspicious. It is not clear how much cheaper shadow-people would be to simulate than real people. It is not even obvious that it is possible for an entity to behave indistinguishably from a real human and yet lack conscious experience. Even if there are such selective simulations, we should not think that we are in one of them unless we think they are much more numerous than complete simulations. There would have to be about 100 billion times as many 'me-simulations' (simulations of the life of only a single mind) as there are ancestor-simulations in order for most simulated persons to be in me-simulations.

There is also the possibility of simulators abridging certain parts of the mental lives of simulated beings and giving them false memories of the sort of experiences that they would typically have had during the omitted interval. If so, one can consider the following (far-fetched) solution to the problem of evil: that there is no suffering in the world and all memories of suffering are illusions. Of course this hypothesis can be seriously entertained only at those times when one is not currently suffering.

Supposing we live in a simulation, what are the implications for us? The foregoing remarks notwithstanding, the implications are not all that radical. Our best guide to how our posthuman creators have chosen to set up our world is the standard empirical study of the universe we see. The revisions to most parts of our belief networks would be rather slight and subtle – in proportion to our lack of confidence in our ability to understand the ways of

posthumans. Properly understood, therefore, the truth of (3) should have no tendency to make us ‘go crazy’ or to prevent us from going about our business and making plans and predictions for tomorrow. The chief empirical importance of (3) at the present time seems to lie in its role in the tripartite conclusion established above.<sup>14</sup> We may hope that (3) is true, since that would decrease the probability of (1), although if computational constraints make it likely that simulators would terminate a simulation before it reaches a posthuman level, then our best hope would be that (2) is true.

If we learn more about posthuman motives and resource-constraints, maybe as a result of developing towards becoming posthumans ourselves, then the hypothesis that we are simulated will come to have a much richer set of empirical implications.

## VII. CONCLUSION

A technologically mature ‘posthuman’ civilization would have enormous computing power. Given this empirical fact, the simulation argument shows that *at least one* of the following propositions is true: (1) the fraction of human-level civilizations that reach a posthuman stage is very close to zero; (2) the fraction of posthuman civilizations that are interested in running ancestor-simulations is very close to zero; (3) the fraction of all people with our kind of experiences who are living in a simulation is very close to one.

If (1) is true, then we will almost certainly go extinct before reaching posthumanity. If (2) is true, then there must be a strong convergence among the courses of advanced civilizations so that virtually none contains any relatively wealthy individuals who desire to run ancestor-simulations and are free to do so. If (3) is true, then we almost certainly live in a simulation. In the dark forest of our current ignorance, it seems sensible to apportion one’s credence roughly evenly between (1), (2), and (3).

Unless we are now living in a simulation, our descendants will almost certainly never run an ancestor-simulation.<sup>15</sup>

*Oxford University*

<sup>14</sup> For some reflections by another author on the consequences of (3), which were sparked by a privately circulated earlier version of this paper, see R. Hanson, ‘How to Live in a Simulation’, *Journal of Evolution and Technology*, 7 (2001).

<sup>15</sup> I am grateful to many people for comments, and especially to Amara Angelica, Robert Bradbury, Milan Ćirković, Robin Hanson, Hal Finney, Robert A. Freitas Jr, John Leslie, Mitch Porter, Keith DeRose, Mike Treder, Mark Walker, Eliezer Yudkowsky, and the anonymous referees.