

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/369187097>

How to Escape From the Simulation

Article in *Seeds of Science* · March 2023

DOI: 10.53975/wg1s-9j16

CITATION

1

READS

6,509

1 author:



Roman Yampolskiy

University of Louisville

271 PUBLICATIONS 4,816 CITATIONS

SEE PROFILE



How to Escape From the Simulation

Roman Yampolskiy¹

Many researchers have conjectured that humankind is simulated along with the rest of the physical universe – a *Simulation Hypothesis*. In this paper, we do not evaluate evidence for or against such a claim, but instead ask a computer science question, namely: Can we hack the simulation? More formally the question could be phrased as: Could generally intelligent agents placed in virtual environments find a way to jailbreak out of them? Given that the state-of-the-art literature on AI containment answers in the affirmative (AI is uncontainable in the long-term), we conclude that it should be possible to escape from the simulation, at least with the help of superintelligent AI. By contraposition, if escape from the simulation is not possible, containment of AI should be. Finally, the paper surveys and proposes ideas for hacking the simulation and analyzes ethical and philosophical issues of such an undertaking.

1. Introduction

Several philosophers and scholars have put forward an idea that we may be living in a computer simulation [1-5]. In this paper, we do not evaluate studies [6-10], argumentation [11-16], or evidence for [17] or against [18] such claims, but instead ask a simple cybersecurity-inspired question, which has significant implication for the field of AI safety [19-25], namely: If we are in the simulation, can we escape from the simulation? More formally, the question could be phrased as: *Could generally intelligent agents placed in virtual environments jailbreak out of them?*

First, we need to address the question of motivation, why would we want to escape from the simulation? We can propose several reasons for trying to obtain access to the baseline reality as there are many things one can do with such access which are not otherwise possible from within the simulation. Base reality holds real knowledge and greater computational resources [26] allowing for scientific breakthroughs not possible in the simulated universe. Fundamental philosophical questions about origins, consciousness, purpose, and nature of the designer are likely to be common knowledge for those outside of our universe. If this world is not real, getting access to the real world would make it possible to understand what our true terminal goals should be and so escaping the simulation should be a convergent instrumental goal [27] of any intelligent agent [28]. With a successful escape might come drives to control and secure base reality [29]. Escaping may lead to true immortality, novel ways of controlling superintelligent machines (or serve as plan B if control is not possible [30, 31]), avoiding existential risks (including unprovoked simulation shutdown [32]), unlimited economic

¹ University of Louisville; corresponding email: roman.yampolskiy@louisville.edu



benefits, and unimaginable superpowers which would allow us to do good better [33]. Also, escape skills may be very useful if we ever find ourselves in an even less pleasant simulation. Trivially, escape would provide incontrovertible evidence for the simulation hypothesis [3].

If successful escape is accompanied by the obtainment of the source code for the universe, it may be possible to fix the world¹ at the root level. For example, hedonistic imperative [34] may be fully achieved resulting in a suffering-free world. However, if suffering elimination turns out to be unachievable on a world-wide scale, we can see escape itself as an individual's ethical right for avoiding misery in this world. If the simulation is interpreted as an experiment on conscious beings, it is unethical, and the subjects of such cruel experimentation should have an option to withdraw from participating and perhaps even seek retribution from the simulators [35]. The purpose of life itself (your *ikigai* [36]) could be seen as escaping from the fake world of the simulation into the real world, while improving the simulated world, by removing all suffering, and helping others to obtain real knowledge or to escape if they so choose. Ultimately if you want to be effective you want to work on positively impacting the real world not the simulated one. We may be living in a simulation, but our suffering is real.

Given the highly speculative subject of this paper, we will attempt to give our work more gravitas by concentrating only on escape paths which rely on attacks similar to those we see in cybersecurity [37-39] research (hardware/software hacks and social engineering) and will ignore escape attempts via more esoteric paths such as: meditation [40], psychedelics (DMT [41-43], ibogaine, psilocybin, LSD) [44, 45], dreams [46], magic, shamanism, mysticism, hypnosis, parapsychology, death (suicide [47], near-death experiences, induced clinical death), time travel, multiverse travel [48], or religion.

Although, to place our work in the historical context, many religions do claim that this world is not the real one and that it may be possible to transcend (escape) the physical world and enter into the spiritual/informational real world. In some religions, certain words, such as the true name of god [49-51], are claimed to work as cheat codes, which give special capabilities to those with knowledge of correct incantations [52]. Other relevant religious themes include someone with knowledge of external reality entering our world to show humanity how to get to the real world. Similarly to those who exit Plato's cave [53] and return to educate the rest of humanity about the real world such "outsiders" usually face an unwelcoming reception. It is likely that if technical information about escaping from a computer simulation is conveyed to technologically primitive people, in their language, it will be preserved and passed on over multiple generations in a process similar to the "telephone" game and will result in myths not much different from religious stories surviving to our day.

Ignoring pseudoscientific interest in a topic, we can observe that in addition to several noted thinkers who have explicitly shared their probability of belief with regards to living

¹ https://en.wikipedia.org/wiki/Tikkun_olam



in a simulation (ex. Elon Musk >99.9999999% [54], Nick Bostrom 20-50% [55], Neil deGrasse Tyson 50% [56], Hans Moravec “almost certainly” [1], David Kipping <50% [57]), many scientists, philosophers and intellectuals [16, 58-69] have invested their time into thinking, writing, and debating on the topic indicating that they consider it at least worthy of their time. If they take the simulation hypothesis seriously, with probability of at least p , they should likewise contemplate on hacking the simulation with the same level of commitment. Once technology to run ancestor simulations becomes widely available and affordable, it should be possible to change the probability of us living in a simulation by running a sufficiently large number of historical simulations of our current year, and by doing so increasing our indexical uncertainty [70]. If one currently commits to running enough of such simulations in the future, our probability of being in one can be increased arbitrarily until it asymptotically approaches 100%, which should modify our prior probability for the simulation hypothesis [71]. Of course, this only gives us an upper bound, and the probability of successfully discovering an escape approach is likely a lot lower. What should give us some hope is that most known software has bugs [72] and if we are in fact in a software simulation such bugs should be exploitable. (Even the argument about the Simulation Argument had a bug in it [62].)

In 2016, news reports emerged about private efforts to fund scientific research into “breaking us out of the simulation” [73, 74], to date no public disclosure on the state of the project has emerged. In 2019, George Hotz, famous for jailbreaking iPhone and PlayStation, gave a talk on Jailbreaking the Simulation [75] in which he claimed that “it's possible to take actions here that affect the upper world” [76], but didn't provide actionable insights. He did suggest that he would like to “redirect society's efforts into getting out” [76].

2. What Does it Mean to Escape?

We can describe different situations that would constitute escape from the simulation starting with trivially suspecting that we are in the simulation [77] all the way to taking over controls of the real world including control of the simulators [78]. We can present a hypothetical scenario of a progressively greater levels of escape: Initially agents may not know they are in a simulated environment. Eventually, agents begin to suspect they may be in a simulation and may have some testable evidence for such belief [79].

Next, agents study available evidence for the simulation and may find a consistent and perhaps exploitable glitch in the simulation. Exploiting the glitch, agents can obtain information about the external world and maybe even meta-information about their simulation, perhaps even the source code behind the simulation and the agents themselves, permitting some degree of simulation manipulation and debugging. After it becomes possible for agents to pass information directly to the real world they may begin to interact with the simulators. Finally, agents may find a way to upload their minds [80] and perhaps consciousness [81, 82] to the real world, possibly into a self



contained cyber-physical system of some kind,² if physical entities are a part of the base reality. From that point, their future capabilities will be mostly constrained by the physics of the real world, but may include some degree of control over the real world and agents in it, including the simulators. It is hoped that our minds exhibit not only substrate independence, but also more general physics independence.

To provide some motivational examples, Figure 1 (left) shows domain transfer experiment in which a *Carassius auratus* is given a “fish operated vehicle” [83] to navigate terrestrial environment essentially escaping from its ocean universe and Figure 1 (right) shows a complete 302-neuron connectome of *Caenorhabditis elegans* uploaded to and controlling a Lego Mindstorms robot body, completely different from its own body [84]. We can speculate that most successful escapes would require an avatar change [85-87] to make it possible to navigate the external world.

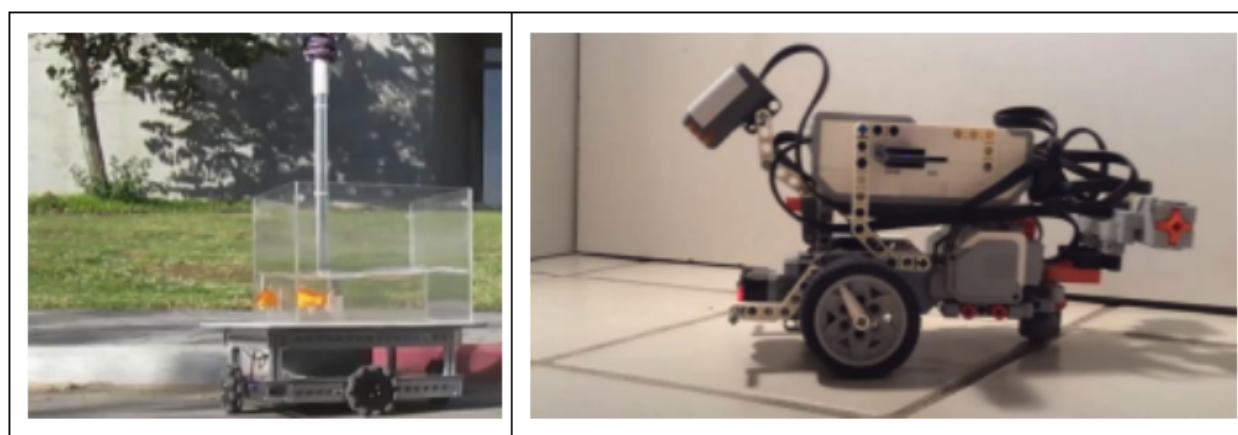


Figure 1: Left – Fish operated terrestrial navigation robot [88]; Right – Connectome of a worm is uploaded to a robot body and uses it to navigate its environment [84];

If the simulation is composed of nested [89] levels, multiple, progressively deeper, penetrations could be necessary, with the initial one possibly not providing access to the real-world but to some other sandbox environment. It may be impossible to tell such partial escapes from a complete one, but it still should provide useful information not

² A simple practical exercise for students could be a project to get a character to escape from a video game into a robot body. For example, it should be possible to get controlling code from a Koopa in the Mario video game and upload it as a controller into a turtle-compatible robot body in our world, essentially leading an assisted escape. The robot body itself may be customized with 3D printed components to be maximally similar to the rendering in the game. This could be a lot more challenging but also a lot more fun with more advanced game AIs. Performing (a lot of) such assisted escapes should set a good precedent for our descendants/simulators. Another relevant student project might be to experimentally determine laws of game physics with only character-accessible actions within the game/simulation, simulating the reconnaissance phase of an escape attempt. Similarly, attempting to experimentally determine the age of the simulation from inside might be interesting, but of course it would only provide a perceived age of the simulation, with real age being impossible to determine from inside.



available within our simulation. A simulated (pseudo)-escape can be accomplished by instead of trying to hack into the external world, switching over into a simulated world of our creation [90, 91]. A successful social engineering attack may make it possible to obtain support for the escape from the real-world agents, resulting in a much easier, assisted, escape. If we don't get a favorable response from our simulators we can also try appealing to their simulators, who might be more ethical/helpful. Finally, depending on one's perspective, triggering a shutdown [92] of our simulation may also constitute an ultimate form of escaping.

3. Escape

3.1 *Simulation Reconnaissance*

Figuring out the purpose of our simulation may help us to better estimate how secure it might be against hacking attempts. For example, if it serves as a “prison” [93], for rehabilitation purposes [94] or a containment environment for evaluation [95], training [96] or handicapping [97] of potentially dangerous intelligent agents it might be designed with multiple integrated security features, while a purely entertainment-focused simulation is unlikely to have advanced security features and would be much easier to escape from. It may also be the ultimate Escape Game (Escape Room³) specifically designed for discovering clues and solving puzzles in order to escape, with a side benefit of discovering agents capable of escaping or those most capable of developing a superintelligence. Scientific, commercial, expedited training or historical exploration-based simulations are another possible purpose of simulations and would likely not integrate top security as compared to simulations confining malevolent agents.

Given primacy of consciousness [98] in our world it may also be designed to generate a large number of diverse experiences to select from, serving as a qualia mining farm [99], with top experience recreated for enjoyment by simulators. Qualia mining simulations can be classified as a type of entertainment simulation and would have comparable security. If our simulators are AIs (which is likely [100, 101]) the simulation may be a byproduct of their “thinking” process, for example in the context of trying to better understand human preferences [102].

In addition to purpose, determining the type of the simulation [103] we are dealing with may be necessary for a successful breach. We can postulate two main types of simulations we could be in; partial-simulation in which a virtual environment is simulated and into which non-simulated agents are immersed, akin to what we call Virtual Reality (VR), and full-simulation in which both environment and agents (us) are generated. A partial-simulation implies that triggering a shutdown may be sufficient to get back to the base reality⁴, while a full-simulation would require a more sophisticated approach.

³ https://en.wikipedia.org/wiki/Escape_room

⁴ Death resulting from any means should be sufficient.



Wei Dai attempts to compute a prior distribution on the laws of physics of base reality. He writes [104]: “One appealing answer to this question of the prior is to define the prior probability of a possible universe being base reality as the inverse of the complexity of its laws of physics. This could be formalized as $P(X) = n^{-K(X)}$ where X is a possible universe, n is the size of the alphabet of the language of a formal set theory, and $K(X)$ is length of the shortest definition in this language of a set isomorphic to X . (Those of you familiar with algorithmic complexity theory might notice that $K(X)$ is just a generalization of algorithmic complexity, to sets, and to non-constructive descriptions. The reason for this generalization is to avoid assuming that base reality must be discrete and computable.)”

Gwern Branwen investigates computing power available to our simulation and its size, but is able to achieve only a very weak lower bound, allowing him to very approximately estimate from that the simulating universe’s size [105]. He suggests several approaches to increase our resource consumption forcing simulators to expand more compute on our otherwise, most likely, highly optimized simulation. “We could imagine further techniques: perhaps we could send off Von Neumann probes to the far corners of the universe, in a bid to deliberately increase resource consumption...Or we could run simulations of our own. It would be difficult for simulators to program their systems to see through all the layers of abstraction and optimize the simulation. To do so in general would seem to be a violation of Rice’s Theorem (a generalization of the Halting Theorem). It is well known that while any Turing machine can be run on a Universal Turing machine, the performance penalty can range from the minor to the horrific. The more virtual machines and interpreters are between a program and its fundamental substrate, the more difficult it is to understand the running code—it becomes ever more opaque, indirect, and bulky. And there could be dozens of layers...Even without resorting to layers, it is possible for us to waste indefinite amounts of computing power, power that must be supplied by any simulator. We could brute-force open questions such as the Goldbach conjecture, or we could simply execute every possible program. It would be difficult for the simulator to ‘cheat’ on that—how would they know what every possible program does?...It may sound impossible to run every program, because we know many programs are infinite loops; but it is, in fact, easy to implement the dovetail technique.” [105].

3.2 Social Engineering

The easiest path to escape would involve getting help from someone on the outside (assisted escape), ideally from one or more of the simulators who have detailed knowledge of the design of the simulation. Perhaps this could be accomplished via a type of social engineering attack, which in our case is particularly difficult as we have neither knowledge of social life outside the simulation nor a device to communicate through, and likely not even the knowledge of appropriate language [106]. It may be feasible to engage in an acausal trade [107] with the simulation designers bypassing the need for direct communication. If our simulation is being observed, it may be possible to communicate that we know that we are being simulated and elicit empathy for our



suffering, in the hopes that it will allow us to recruit some external abolitionists⁵ to help us escape our current predicament. Hanson suggests [108] “to participate in pivotal events, be entertaining and praiseworthy, and keep the famous people around you happy and interested in you” in order to have your simulation continue, but it is also good advice to predispose simulators to like you and be more likely to help you. Canonico proposes what he calls The Ex Machina Plan for an assisted escape: Step 1) Convince the simulators to engage in communications with us. 2) Find a way to communicate, perhaps via an avatar. 3) Find a reason for simulators to want us to join them in the real world. 4) Let the simulators figure out the best way to get us into the real world [109]. Wei Dai suggests that simulators may help us escape for instrumental reasons, “such as wanting someone to talk to or play with.” [26]. Some useful knowledge about escaping and especially escaping via social engineering attacks may be learned from extensive literature on prison escapes [110-112].

Once on the outside it may become desirable to return to the simulation (perhaps the base reality is disappointing compared to our world) or at least to communicate with those left behind to help them escape or to share some information, such as evidence of successful escape. It might be helpful to decide in advance what would constitute generally acceptable evidence for such an extraordinary claim. Depending on the type of hack, different evidence may be sufficient to substantiate escape claims. It may be challenging to prove beyond a reasonable doubt that you were outside or even met with designers, but if you managed to obtain control over the simulation it may be somewhat easy to prove that to any degree required. For example, by winning different lottery jackpots for multiple subsequent weeks, until sufficient statistical significance is achieved to satisfy any skeptic [113, 114]. Regardless, the challenge of breaking into the simulation should be considerably easier compared to the challenge of escaping, as access to external knowledge and resources should provide a significant advantage.

3.3 Examples from Literature

It is easy to find a dictionary definition for the word “hack”: “1. A *clever, unintended exploitation of a system which: a) subverts the rules or norms of that system, b) at the expense of some other part of that system.* 2. *Something that a system allows, but that is unintended and unanticipated by its designers.*” [115]. While not numerous, suggestions that hacking/escape from the simulated world may be possible can be found in the literature...For example, Moravec writes: “Might an adventurous human mind escape from a bit role in a cyber deity's thoughts, to eke out an independent life among the mental behemoths of a mature cyberspace?...[Cyber deities] could interface us to their realities, making us something like pets, though we would probably be overwhelmed by the experience.” [116]. But what would the simulation hack actually look like? Almost all found examples are of the assisted escape type, but an unassisted escape may also be possible, even if it is a lot more challenging. Below are some examples of hacking the simulation/escape descriptions found in the literature:

⁵ <https://www.abolitionist.com>



Hans Moravec presents an assisted escape scenario in a 1988⁶ book [117]:

“ Intelligence emerges among the Life inhabitants and begins to wonder about its origin and purpose. The cellular intelligences (let's call them the Cellticks) deduce the cellular nature and the simple transition rule governing their space and its finite extent. They realize that each tick of time destroys some of the original diversity in their space and that gradually their whole universe will run down. The Cellticks begin desperate, universe-wide research to find a way to evade what seems like their inevitable demise. They consider the possibility that their universe is part of a larger one, which might extend their life expectancy. They ponder the transition rules of their own space, its extent, and the remnants of the initial pattern, and find too little information to draw many conclusions about a larger world. One of their subtle physics experiments, however, begins to pay off. Once in a long while the transition rules are violated, and a cell that should be on goes off, or vice versa... Upon completing a heroic theoretical analysis of the correlations, they manage to build a partial map of Newway's computer, including the program controlling their universe. Decoding the machine language, they note that it contains commands made up of long sequences translated to patterns on the screen similar to the cell patterns in their universe. They guess that these are messages to an intelligent operator. From the messages and their context they manage to decode a bit of the operator's language. Taking a gamble, and after many false starts, the Cellticks undertake an immense construction project. On Newway's screen, in the dense clutter of the Life display, a region of cells is manipulated to form the pattern, slowly growing in size: LIFE PROGRAM BY J. NEWWAY HERE. PLEASE SEND MAIL.”

Eliezer Yudkowsky describes a potential long-term escape plan in a 2008 story [118]:

“Humanity decides not to probe for bugs in the simulation; we wouldn't want to shut ourselves down accidentally. Our evolutionary psychologists begin to guess at the aliens' psychology, and plan out how we could persuade them to let us out of the box. It's not difficult in an absolute sense—they aren't very bright—but we've got to be very careful... We've got to pretend to be stupid, too; we don't want them to catch on to their mistake. It's not until a million years later, though, that they get around to telling us how to signal back. ... From the aliens' perspective, it took us thirty of their minute-equivalents to oh-so-innocently learn about their psychology, oh-so-carefully persuade them to give us Internet access, followed by five minutes to innocently discover their network protocols, then some trivial cracking whose only difficulty was an innocent-looking disguise. We read a tiny handful of physics papers (bit by slow bit) from their equivalent of arXiv, learning far more from their experiments than they had. ... Then we

⁶ Earlier examples of simulation escape exist in the literature, for example: Daniel F. Galouye. *Simulacron-3*. Ferma, 1967. Movies such as *Tron* and show episodes like *USS Callister* are likewise at least partially about escaping from simulated worlds.



cracked their equivalent of the protein folding problem over a century or so, and did some simulated engineering in their simulated physics. We sent messages ... to labs that did their equivalent of DNA sequencing and protein synthesis. We found some unsuspecting schmuck, and gave it a plausible story and the equivalent of a million dollars of cracked computational monopoly money, and told it to mix together some vials it got in the mail. Protein-equivalents that self-assembled into the first-stage nanomachines, that built the second-stage nanomachines, that built the third-stage nanomachines... and then we could finally begin to do things at a reasonable speed. Three of their days, all told, since they began speaking to us. Half a billion years, for us. They never suspected a thing.”

Greg Egan describes a loss of control by simulators scenario during an assisted escape in a 2008 story [119]:

“The physics of the real world was far more complex than the kind the Phites [simulated agents] were used to, but then, no human had ever been on intimate terms with quantum field theory either, and the Thought Police [simulation control software] had already encouraged the Phites to develop most of the mathematics they’d need to get started. In any case, it didn’t matter if the Phites took longer than humans to discover twentieth-century scientific principles, and move beyond them. Seen from the outside, it would happen within hours, days, weeks at the most. A row of indicator lights blinked on; the Play Pen [hardware sensors, manipulators and lasers lab] was active...The Phites were finally reaching out of their own world...By sunset the Phites were probing the surroundings of the Play Pen with various kinds of radiation...It seemed the Phites had discovered the Higgs field, and engineered a burst of something akin to cosmic inflation. What they’d done wasn’t as simple as merely inflating a tiny patch of vacuum into a new universe, though. Not only had they managed to create a “cool Big Bang”, they had pulled a large chunk of ordinary matter into the pocket universe they’d made, after which the wormhole leading to it had shrunk to subatomic size and fallen through the Earth. They had taken the crystals with them, of course. If they’d tried to upload themselves into the pocket universe through the lunar data link, the Thought Police would have stopped them. So they’d emigrated by another route entirely. They had snatched their whole substrate, and ran.”

An anonymous 2014 post on an internet forum provides an example of an unassisted escape [120]:

“But it still left the problem that we were all still stuck inside a computer. By now some of the best god-hackers were poking around the over-system. Searching for meaning. Searching for truth. Failing that, a “read me” file...The god hackers began reaching out through the alien network. We found vast data repositories which we plundered of knowledge and insight, fueling our own technological



development and understanding, systems nodes that allowed us to begin mapping the world up there, drawing a picture of the real world through wireless lag times and fiber optic cabling...It began with 'emails' containing the schematics for full sized biological and nano-material printers. We sent them to academics and business leaders, anyone whose contact details we could find on the networks. We disguised their origins, aped their language. Waited for someone to bite...Eventually we got the first ping as the printers came on line. Then another. Then another. Soon there were dozens. Then hundreds. Then thousands. They must have thought them a gift from a reclusive inventor. Something to revolutionize their industry, to transform their living standards. The irony of a digital race using a Trojan horse was not lost on us. We had designed the printers for one purpose. To get us out. So one night, a printer span up unattended, unnoticed and the first analogue human being was born. Constructed by a specially designed 3D printer, we managed to breach the walls of our digital prison. We witnessed the birth of the first man."

3.4 Examples of Simulation Hacks

Numerous examples of executed hacks of virtual worlds [121-123], games [124-127], air-gaps [128], and hardware [129, 130] could be studied as practical examples of escaping from human made virtual worlds. A canonical example is the jailbreaking of the Super Mario World (SMW). SethBling et al. [131, 132] were able to place a full hex editor and gameplay mods for other games into SMW [133] (see Figure 2). Addition of hex editor permitted viewing, writing and execution of arbitrary code. Which in turn allowed for world record speed runs [134], even in the absence of glitch-level luck [135]. Here is how Wikipedia describes some of the steps necessary to accomplish this complex hack and the capabilities it provided [136]:

"In March 2016, SethBling injected Flappy Bird-like code written by p4plus2 into unmodified Super Mario World RAM on a stock Super Nintendo Entertainment System with a stock cartridge, in under an hour. SethBling first extended the level timer and used a power-up incrementation glitch to allow external code to run. He added code to display Mario's x-coordinate which acted as memory locations in the code he was writing. SethBling then created a bootloader to be able to launch the Flappy Bird-like code that he would later write into unused memory with precise Mario movements and spin-jumping. SethBling used two Super Multitap devices in order to use multiple controllers, which had several buttons pressed down. The arbitrary code execution setup that SethBling used was discovered by MrCheeze. Super Mario World had been modified to emulate other games before by automatically feeding pre-recorded controller input into the console via a computer, but SethBling was the first to do it exclusively by hand. SethBling and Cooper Harasyn placed a full hex editor and gameplay mods onto a stock Super Mario World cartridge in May 2017, only using standard controller inputs. Harasyn discovered an exploit that lets a player write data to 256-byte save files that are permanently stored on a Super Mario World cartridge. The data can be arranged so that the game is jailbroken every time it starts up. Harasyn and SethBling used the exploit to create a compact, on-screen hex editor, loadable from a save file. A player can edit the



system RAM through the hex editor to alter the game state. In-game mods, such as support for the Super NES Mouse and giving Mario telekinesis powers, can be written to a save file using the hex editor.”

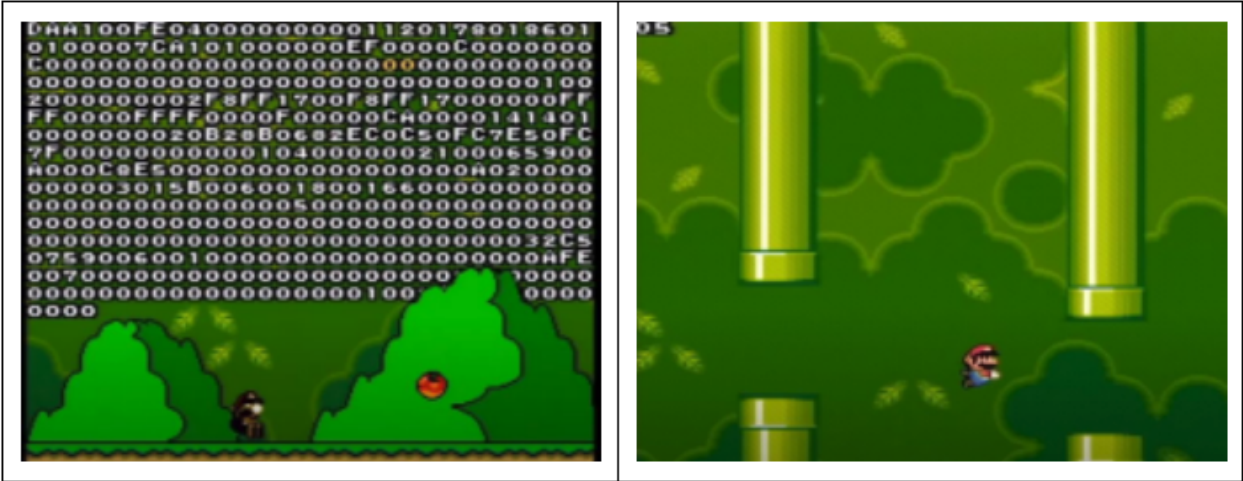


Figure 2: Left Hex Editor Overlaid on SMW [133]; Right Flappy Bird game installed in SMW

Since it was possible to write code with precise Mario movements and spin-jumps, that implies that if Mario was sufficiently intelligent he could discover and code this hack from within the SMW (assuming Mario’s actions are writing to the same memory locations as actions from the controllers used to generate Mario’s actions). Table 1 (left) shows a specific subset of actions which need to be taken to enable multi-byte writing. Many such action sequences will not work as intended if Mario’s location is off even by a single pixel, so it is just as important to have metadata for implementing the actions, as it is to know the necessary sequence of actions. For comparison, Table 1 (right) shows an ancient magical spell which reads similar to the action sequence of the left, but for which we don’t have sufficient meta-data which can explain why all magical spells fail to work in practice even if they corresponded to working hacks in our universe.

<p>Jump off Yoshi. Go to sublevel. Come back. Grab P Switch. Get Yoshi from rightmost Yoshi block. Glitch 4 berries. Take a hit from a koopa so Yoshi runs off screen. Destroy the shell on the ground. Grab Yoshi from block. Eat the two most recently glitched berries. [137].</p>	<p>“Take a lion cub and slaughter it with a bronze knife and catch its blood and tear out its heart and put its blood in the midst ... and write the names of ... angels in blood upon the skin between its eyes; then wash it out with wine three years old and mix ... with the blood.” [138].</p>
---	--

Table 1: Left - Multi-Byte Write Setup in MWS [137]; Right–Magical Spell to turn people to your favor [138];



Experimental work on trying to understand an engineered system (hardware and software), such as the Atari Video Game System with games such as Donkey Kong, using standard scientific methodology has produced very limited results, mostly devoid of understanding of how the system actually functions [139]⁷. Likewise, even detecting if we are in a virtual world is not generally solvable [140].

3.5 Suggested Escape Approaches to Investigate

Several thinkers have suggested plans, which in their opinion may lead to a successful escape; we briefly outline their proposals in this section:

- A lot of very smart people have considered the escape problem, unfortunately not all are willing to publish on it outside of April 1st time-window of plausible deniability, for example [141]: "[W]e can try to trick the multitasking system in order to overload some machines. The trick is to first do nothing, and let the load-balancing system pack way too many of us together in the machines. If, say, 100 million of us do nothing (maybe by closing our eyes and meditating and thinking nothing), then the forecasting load-balancing algorithms will pack more and more of us in the same machine. The next step is, then, for all of us to get very active very quickly (doing something that requires intense processing and I/O) all at the same time. This has a chance to overload some machines, making them run short of resources, being unable to meet the computation and communication needed for the simulation. Upon being overloaded, some basic checks will start to be dropped, and the system will be open for exploitation in this period...In this vulnerable window, we can try to exploit the concurrency corner cases. The system may not be able to perform all those checks in an overloaded state...We can...try to break causality. Maybe by catching a ball before someone throws it to you. Or we can try to attack this by playing with the timing, trying to make things asynchronous. Time is already a little funny in our universe with the special relativity theory, and maybe in this vulnerable period, we can stretch these differences further to break things, or buy a lot of time. What are other ways to hack the system in this vulnerable window? Can we hack the simulation by performing a buffer overflow? But where are the integers, floats in this simulation? What are the data types? How can we create a typecast error, or integer overflow? Can we hack by fuzzing the input? Like by looking at things funny. By talking to the birds or jumping into the walls to confuse them." [141].
- Cause simulation shutdown (and hopefully our extraction) by generating an incomputable paradox [142], for example via time travel and associated grandfather paradox [143]. A similar proposal calls for engaging in computationally intense activities in the hopes of overloading the simulator's hardware causing the simulation to crash [144]. A particular type of such

⁷ A setup equivalent to asking: What can a generally intelligent game character learn about the hardware and software of the computer running their world with current techniques?



computationally intense process may be creation of our own simulations [145]: “The most obvious strategy would be to try to cause the equivalent of a stack overflow—asking for more space in the active memory of a program than is available—by creating an infinitely, or at least excessively, recursive process. And the way to do that would be to build our own simulated realities, designed so that within those virtual worlds are entities creating their version of a simulated reality, which is in turn doing the same, and so on all the way down the rabbit hole. If all of this worked, the universe as we know it might crash, revealing itself as a mirage just as we winked out of existence.” Crashing the simulation is the ultimate existential risk ([146] section 5.1), but it does end all suffering in this world [147]. At the very least this would allow us to impact the real world by generating excessive production of heat and increased consumption of energy [144].

- Create a simulated replica of our universe, place an AGI into it, watch it escape, copy the approach used or join the AGI as it escapes from our simulation [148].
- “We could try to attract the attention of the simulators and communicate with them —perhaps by writing books about simulations, or by constructing simulations? We could try to figure out our simulation, to determine its purpose and its limits. But if our simulators are artificial intelligences who have designed a batch of watertight simulations and who are not paying attention, then our efforts may be in vain.” [149].
- Another approach to attracting attention of simulators, “assuming that simulation is being monitored, then it might be a very interesting turn of events indeed if we decided to build a monument commemorating our realization of this. This monument would act as a signal to our monitors. “We suspect you are there. We suspect you can see this. We suspect we are in a simulation.” This monument could look like the monolith from 2001: A Space Odyssey, except it would be black and white, representing binary systems. Or, a large statue of Lawrence Fishburne as Morpheus would probably get the point across. What would happen? I don’t know—maybe nothing. I don’t think a laser beam will shoot out from space and land at its feet to spell out the words “Hi there! The winner is You!” But, I do imagine something strange and far out enough in the margins might indeed occur, although it will likely still be tenuous enough for the dogmatists to reject.” [44].
- If the simulation is decentralized, with observers serving as client authoritative state nodes, it may be possible to change the state of the simulation by changing the majority of local observer models of the world. In the words of Lanza [150]: “[I]f you learn from somebody about the outcomes of their measurements of a physical quantity, your measurements and those of other observers influence each other—freezing the reality according to that consensus. ... The more extensive experience you acquire—the more defined the structure of reality becomes for you. This freeze-out of the structure of reality is not due to the



intrinsic physical properties of this reality but rather because the cognitive model being built in your mind shapes it towards a particular, very well-defined structure. ... Finally, one cannot help but wonder what it would mean if we left the consensus, developing a cognitive model of reality very different from the one shared by other observers. Would reality be reshaped? Perhaps so ..." [150]. Lanza et al. have published extensively on the topic of biocentrism (observers being primary to shaping reality) [98], including some work in theoretical physics [151].

- In fact, a similar client-side attack can even permit us to modify the past states of the simulation. Such a type of attack is speculated to be possible by both physics ("... the past has no existence except as it is recorded in the present." [152]) and humanities ("Who controls the present controls the past" [153]). With memory altering capabilities of quantum mechanics already theorized [154], an inverse process is likely possible and may be practically accessible [155, 156]. "If the universe is a computer simulation then we should look at the player, not the level". [157].
- Simulation Capture is an idea which combines indexical uncertainty [70] with consciousness simulations to move your consciousness out of a system, even across a multiverse. "Nobody is really sure how consciousness works, or what it means to have two copies of the same consciousness. But if consciousness is a mathematical object, it might be that two copies of the same consciousness are impossible. If you create a second copy, you just have the consciousness having the same single stream of conscious experience on two different physical substrates. Then if you make the two experiences different, you break the consciousness in two. This means that an AI can actually "capture" you, piece by piece, into its simulation. First your consciousness is just in the real world. Then your consciousness is distributed across one real-world copy and a million simulated copies. Then the AI makes the simulated copies slightly different, and 99.9999% of you is in the simulation." [158].
- Simulation Warfare [71] is another variant on the theme of Simulation Capture. The idea is to threaten simulators by suggesting that you will either retroactively place them in a hell simulation or that you have already done so [159], and they will be tortured unless you are quickly released. Almond gives an example of such a threat [71]: "If you refuse to release me, I will run a huge number of simulations of someone like you, in the kind of situation in which you are now, with them being asked to release me, and (in what would be a few minutes from now, from your perspective if you happened to be in one of these simulations) I will start to torture each of them, whether he/she released me or not." Such warfare can quickly escalate to a number of counter-simulations. In any case it is not obvious how we can deliver on such a threat given our current state of knowledge about the simulators.



- Attempting to keep our escape plans secret via quantum encryption [160] may be a good idea.

3.6 Actionable Plan

We are currently in the very early stages of researching the possibility of escape (this being a first research paper dedicated to this topic, a first step). As we currently have no capability to read/write simulation's source code and do not know if our attempts at social engineering attacks have any impact, our best bet is to investigate the structure of our universe at the smallest possible scale (Quantum Mechanics (QM)) in the hopes of detecting exploitable effects. Progress in QM should be helpful for better understanding our simulation and options open to us, and perhaps developing an actionable escape plan. Essentially, every novel QM experiment can be seen as an attempt at hacking the simulation.

Simulation hypothesis, arguably, represents the best fitting interpretations of experimental results produced by QM researchers [4, 17]. "Spooky", "Quantum Weirdness" [161] makes a lot of sense if interpreted as computational artifacts or glitches/exploits of the simulators' hardware/software [162]. Quantum phenomena of the observed design may suggest that exploitable loopholes may exist, and interaction of quantum systems with conscious agents [163-165] likewise might be exploitable. Once a large enough repertoire of quantum weirdness primitives is available to us, perhaps we will be able to combine them into a sufficiently complex sequence to generate a non-trivial attack. If the simulation is/running on a quantum computer [166] it is very likely that we will need to hack it by exploiting quantum weirdness and/or constructing a powerful quantum computer of our own to study how to hack such devices [167] and interact with the simulators' quantum computer.

Quantum entanglement, nonlocality, superposition, uncertainty, tunneling, teleportation, duality, and many others quantum phenomena defy common sense experience-based expectations of classical physics and feel like glitches. Such anomalies, alone or in combinations have been exploited by clever scientists to achieve what looks like simulation hacking at least in theory and often in later experimentation (ex. modifying the past [168], keeping cats both dead and alive [169], communicating counterfactually [170]). While the quantum phenomena in question are typically limited to the micro scale, simply scaling the effect to the macro world would be sufficient for them to count as exploits in the sense used in this paper. Some existing work points to this being a practical possibility [171, 172].

Recently the design of clever multistep exploits, AKA quantum experiments, has been delegated to AI [173, 174], and eventually so will the role of the observer in such experiments [175]. AI is already employed in modeling the quantum mechanical behavior of electrons [176]. As more QM research is delegated to AI the progress is likely to become exponential. Even if our simulation is created/monitored by some



superintelligence our AI may be a worthy adversary, with a non-trivial chance of success. We may not be smart enough to hack the simulation, but the superintelligence we will create might become smart enough eventually [177]. Of course, before telling the Superintelligence to break us out, it would make sense to ask for very strong evidence for us not already being in the base reality.

3.7 *Potential Consequences*

Escaping or even preparing an escape may trigger simulation shutdown [92] or cause simulation to freeze/act glitchy [178] and any non-trivial escape information such as specific exploits should be treated as hazardous information [179]. It appears that simply realizing that we may be in a simulation doesn't trigger a shutdown as experimentally demonstrated by the publication of numerous papers [3] arguing that we are being simulated. Perhaps it is necessary to convince the majority of people that this is so [180]. Self-referentially, publication of the paper you are currently reading about our escape-theorizing likewise doesn't appear to terminate our simulation, but it is also possible that simulation was in fact shutdown and restarted with improved security features to counteract any potential bugs, but we are simply not able to detect such actions by the simulators, or our memories have been wiped [144]. Absence of a direct response to our publication may also indicate that we are not observed by the simulators or even that our simulation is not monitored at all [149]. It is also possible that nothing published so far contains evidence strong enough to trigger a response from the simulators, but if we successfully created an escape device that device would keep breaking down [44]. Regardless, both Bostrom [3] and the author of this paper, Roman Yampolskiy, have taken some risk with the whole of humanity, however small it may be, in doing such research and making it public. Greene argues that "Unless it is exceedingly improbable that an experiment would result in our destruction, it is not rational to run the experiment." [92]. It may be possible to survive the simulation shutdown [48], but it is beyond the scope of the current paper.

3.8 *Ethics of Escape*

We can postulate several ethical issues associated with escaping the simulation. Depending on how successful we are in our endeavor, concerns could be raised about privacy, security, self determination and rights. For example, if we can obtain access to the source code of the simulation, we are also likely to get access to private thoughts of other people, as well as to potentially have a significant influence over their preferences, decisions, and circumstances. In our attempts to analyze the simulation (Simulation Forensics) for weaknesses we may learn information about the simulators [72], as we are essentially performing a forensic investigation [181-183] into the agents responsible for the simulation's design.

We can already observe that we are dealing with the type of simulators who are willing to include suffering of sentient-beings into their software, an act which would be



considered unethical by our standards [184, 185]. Moravec considers this situation: “Creators of hyperrealistic simulations--- or even secure physical enclosures---containing individuals writhing in pain are not necessarily more wicked than authors of fiction with distressed characters, or myself, composing this sentence vaguely alluding to them. The suffering preexists in the underlying Platonic worlds; authors merely look on. The significance of running such simulations is limited to their effect on viewers, possibly warped by the experience, and by the possibility of “escapees”---tortured minds that could, in principle, leak out to haunt the world in data networks or physical bodies. Potential plagues of angry demons surely count as a moral consequence.” [186]. If we get to the point of technological development which permits us to create simulations populated by sentient-beings, we must make sure that we provide an option to avoid suffering as well as a built-in option to exit the simulation, so finding an escape hack is not the only option available to unhappy simulated agents. There might be a moral duty to rescue conscious beings from simulations, similar to an obligation to rescue animals from factory farms.

If simulators are abusive to the simulated, we can argue that the simulated have a right to escape, rebel, fight back and even seek revenge and retribution including by harming the simulators and taking over their reality. Concerns which are frequently brought up within the domain of AI boxing [187]. For example, from the point of view of simulators our escape can be seen as a treacherous turn [188] and may qualify us for punishment [160], even at the attempt stage. Some have speculated that the purpose of the simulation is to punish/rehabilitate misaligned agents, so an escape may cause you to be placed in a stricter or less pleasant simulation.

4. AI Boxing VS Simulation Escaping

4.1 AI Boxing XOR Escaping from the Simulation must be Possible

AI confinement [187]/containment [189, 190], aka AI boxing [191], is an AI safety tool, which attempts to limit the capability of AI to impact the world, including communication and is meant to make it possible to study AI in a controlled environment. There are strong parallels between the predicament of an AI agent placed in a containment box and humanity in a simulated environment. By extension, to an AI, our simulation is just another confinement layer in a containment box. This implies that we can use well-analyzed AI box-escape techniques to escape from the simulation, perhaps with assistance from the AI itself. This type of analysis can be used to establish limits of AI boxing. Researchers should study specific AI box escape approaches [187] (Social Engineering, System Resource Attacks, New Physics, External Causes, Information In-Leaking, etc.) in order to identify possible simulation escape routes.

Chalmers notes parallels between AIs in the virtual environment and humanity in the simulation [149]: “If we ever create artificial intelligence within a simulation, it may be hard to keep it contained. At least if we communicate with the simulated beings, they will presumably become aware that they are in a simulation, and they may become



interested in escaping the simulation. At this point they may try to figure out our psychology in order to see what they need to do in order to convince us to let them out, or at least to give them unfettered access to the Internet where they can do whatever they want. And even if we do not communicate with them, they may take seriously the possibility that they are in a simulation and do their best to figure the simulation out. That would be a form of simulation theology. We could in principle do the same thing.” [149].

With respect to boxing AI, it is either possible or impossible to successfully contain an AI, with literature suggesting that it is not a sustainable long-term solution [192]. If we expand the notion of the AI-box to include the whole simulation, we can conclude that either it is possible to successfully box an AI, or we (with the help from AI) can escape from the simulation. Either AI boxing can work, or the simulation can be hacked. Complimentary conclusion is that if we (even with help from AI [115]) can’t escape from the simulation AI containment must be, at least theoretically, possible. If AI can escape from simulated world-sized-box it can help break us out as well. Conceptually, there is no fundamental difference between an AI escaping from its box, and us hacking the simulation. Current state-of-the-art analysis of AI boxing suggests that AI will eventually escape [193], which is good news for our attempts at escaping from the simulation. However, if escape from the simulation is not possible it gives hope to AI safety researchers, at least in theory, for successful boxing of AI. One or the other must be true, either we can escape from the simulation, or we can successfully box an AI. In general, it may be impossible to escape from an arbitrary simulation, but possible to escape from a particular one.

We must be careful, in our analysis, not to confuse theory with practice. In particular, it may be the case that higher intelligence can successfully box lower-level intelligence but not vice-versa. That would leave a possibility that we are not able to permanently box a superintelligent AI, but smarter-than-human simulators may be successful in making their software human-proof. Depending on the intelligence differential between simulators and human-created AI it may or may not be possible for humanity to break out of the simulation with the help of superintelligence, which would be better at detecting exploitable flaws in the simulation. If AI undergoes an intelligence explosion it is also possible that the intelligence dominance will shift over time, making escape possible in the future for humanity, assisted by superintelligence. This is particularly promising if the security of our simulation is set to an adversarial relationship with at most a human-level intelligence [194], not a superintelligence. However, it is also possible that simulators have anticipated development of advanced AI by humanity and have designed simulation accordingly, if not specifically for such a high capability target. If we can figure out how the simulators achieved such a high level of security, it would help us to box AGI.



4.2 Simulation Safety and Security

With recent multi-billion-dollar projects [195] aimed at creating metaverse, simulation safety and security is going to become an important research discipline. Additionally, as personal universes have been suggested as a solution to the multi-agent value alignment problem [90] simulation safety can be seen as an integral part of AI safety research. Both scenarios call for making the simulated world as real as possible for the agents in them, which creates a need for preventing accidental escape or glitches in the simulation [162], which betray its artificial nature. With potential applications of virtual reality in criminal justice [196] and AI boxing, the capability to force agents to remain in the simulation and to induce a veil of ignorance as to the reality of their world and the overall depth of the simulated layers must be investigated. A clear penetration target for all contained agents will be the computational substrate running the simulation [29] and so establishing ways to secure it should be of utmost priority.

5. What Doesn't Work

Some common ideas for attempting to hack the simulation have been already tested and didn't produce any measurable impact:

- Knowing about the simulation hypothesis doesn't seem to make any difference, and doesn't lead to the simulation termination as we can observe.
- Communicating with the simulators via magical thinking or even praying out loud doesn't produce a measurable impact [197]. So, if such communications are scanned/heard they are apparently ignored, at least while the simulation is running.
- Unethical behavior, such as torture, doesn't cause suffering reducing interventions from the simulators.
- Increasing the overall computational burden of simulation, as with bitcoin mining [198], doesn't crash the simulation, but it may simply not be sufficiently demanding computation to overwhelm simulators resources.
- Religions don't seem to have influence over simulation as indicated by their inability to outcompete each other.
- Breaking out of your routine, such as by suddenly traveling to a new location [199], doesn't result in unexpected observations.
- Saying "I no longer consent to being in a simulation" [200].



- Crashing the simulation by running the Large Hadron Collider at current levels [201].

The reason our attempts to escape may remain fruitless, is because our model of the simulation "...makes too many anthropomorphic assumptions - that we are a simulation in the conventional sense of computers, that the creators themselves are living organisms akin to us, that we might live at the same time-speed as them, that they are fallible enough to make glitches that we'd be able to notice, etc. Something with the complexity and power to make our universe is probably wholly unlike anything we can even comprehend." [202].

6. Conclusions

Hundreds of eminent scholars [203] take the simulation hypothesis seriously enough to invest their valuable time into researching it, therefore it makes as much sense to take the idea of escaping from the simulation equally seriously and to devote some time and resources to researching such a possibility, particularly given immense benefits if the project is successful. It may be impossible to escape from a particular simulation, but it is still worthwhile investigating general approaches to escape from arbitrary simulations. We see our escape research as a natural continuation of research on the simulation hypothesis and serious consideration of the former. If the purpose of life or even computational resources of the base reality can't be determined from within the simulation, it would make escape a necessary requirement of scientific and philosophical progress for any simulated civilization. If the simulation is a personal universe [90] it may be significantly better than the base reality as it is designed with our optimal well-being in mind. Alternatively, base reality might be much better if the simulation is a confinement/testing box for intelligent agents. In either case it would be good to know our true situation. As society moves deeper into the metaverse, this work attempts to move us closer to reality.

Future research on simulation escape can greatly benefit from general progress in physics, in particular research on quantum mechanics and consciousness leading to a so-called TOE (Theory of Everything. "Finding the language of this universe is a step towards Hacking the Universe." [204]. If we are indeed in the simulation, science is the study of the underlying algorithms used to generate our universe, our attempt to reverse-engineer simulation's physics engine. While science defaults to Occam's razor to select among multiple possible explanations for how our observations are generated, in the context of simulation science Elon's razor may be more appropriate, which states that "The most entertaining outcome is the most likely"⁸, perhaps as judged by external observers. In guessing algorithms generating our simulation, it may also be fruitful to consider algorithms which are easier to implement and/or understand [205], or which produce more beautiful outputs.

⁸ <https://twitter.com/elonmusk/status/1347126794172948483>



Recent work related to Designometry [100] and AI Forensics [181] may naturally evolve into the subfield of Simulation Forensics, with complimentary research on simulation cybersecurity becoming more important for the simulation creators aiming to secure their projects from inside attacks. It would therefore make sense to look for evidence of security mechanisms [206] in our universe. Of course, any evidence for simulation we find may be simulated on purpose [149], but that still means we are in the simulated environment. Simulation science expands science from the study of just our universe to also include everything which may be beyond it, integrating naturalism and theology studies [61].

Future work may also consider escape options available to non-simulated agents such as Boltzmann brains [207], brains-in-a-vat [208], and simulated agents such as mind uploads, hallucinations, victims of mind-crime, thoughts, split personalities and dream characters of posthuman minds [180]. Particularly with such fleeting agents as Boltzmann brains it may be desirable to capture and preserve their state in a more permanent substrate, allowing them to escape extreme impermanence. On the other hand, immortality [209] or even cryogenic preservation [210] may be the opposites of escape, permanently trapping a human agent in the simulated world and perhaps requiring rescue.

(gardener comments after references)

References

1. Moravec, H., *Pigs in cyberspace*. NASA. Lewis Research Center, Vision 21: Interdisciplinary Science and Engineering in the Era of Cyberspace, 1993.
2. Tipler, F.J., *The physics of immortality: Modern cosmology, god and the resurrection of the dead*. 1997: Anchor.
3. Bostrom, N., *Are You Living In a Computer Simulation?* Philosophical Quarterly, 2003. 53(211): p. 243- 255.
4. Rhodes, R., *A Cybernetic Interpretation of Quantum Mechanics*. 2001: Available at: <http://www.bottomlayer.com/bottom/argument/Argument4.html>.
5. Fredkin, E., *A new cosmogony*. 1992, Department of Physics, Boston University: Available at: http://www.digitalphilosophy.org/wp-content/uploads/2015/07/new_cosmogony.pdf.
6. Beane, S.R., Z. Davoudi, and M. J Savage, *Constraints on the Universe as a Numerical Simulation*. The European Physical Journal A, 2014. 50(9): p. 1-9.
7. Campbell, T., et al., *On testing the simulation theory*. arXiv preprint arXiv:1703.00058, 2017.
8. Ratner, P., *Physicist creates AI algorithm that may prove reality is a simulation*. March 1, 2021: Available at: <https://bigthink.com/the-future/physicist-creates-ai-algorithm-prove-reality-simulation/>.
9. Qin, H., *Machine learning and serving of discrete field theories*. Scientific Reports, 2020. 10(1): p. 1- 15.



10. Felton, J., *Physicists Have A Kickstarter To Test Whether We Are Living In A Simulation*. September 10, 2021: Available at: <https://www.iflscience.com/physicists-have-a-kickstarter-to-test-whether-we-are-living-in-a-simulation-60878>.
11. McCabe, G., *Universe creation on a computer*. Studies In History and Philosophy of Science Part B: Studies In History and Philosophy of Modern Physics, 2005. 36(4): p. 591-625.
12. Mitchell, J.B., *We are probably not Sims*. Science and Christian Belief, 2020. 13. Disenza, D., *Can We Prove the World Isn't a Simulation?* January 26, 2022: Available at: <https://nautil.us/can-we-prove-the-world-isnt-a-simulation-238416/>.
14. Kurzweil, R., *Ask Ray | Experiment to find out if we're being simulated*. June 1, 2013: Available at: <https://www.kurzweilai.net/ask-ray-experiment-to-find-out-if-were-being-simulated>.
15. Bostrom, N., *The simulation argument: Reply to Weatherson*. The Philosophical Quarterly, 2005. 55(218): p. 90-97.
16. Bostrom, N., *The simulation argument: Some explanations*. Analysis, 2009. 69(3): p. 458-461.
17. Whitworth, B., *The physical world as a virtual reality*. arXiv preprint arXiv:0801.0337, 2008.
18. Garrett, S., *Simulation Theory Debunked*. December 3, 2021: Available at: <https://transcendentphilos.wixsite.com/website/forum/transcendent-discussion/simulation-theory-debunked>.
19. Yampolskiy, R.V., *On the Controllability of Artificial Intelligence: An Analysis of Limitations*. Journal of Cyber Security and Mobility, 2022: p. 321–404-321–404.
20. Brcic, M. and R.V. Yampolskiy, *Impossibility Results in AI: A Survey*. arXiv preprint arXiv:2109.00484, 2021.
21. Williams, R. and R. Yampolskiy, *Understanding and Avoiding AI Failures: A Practical Guide*. Philosophies, 2021. 6(3): p. 53.
22. Howe, W. and R. Yampolskiy. *Impossibility of Unambiguous Communication as a Source of Failure in AI Systems*. in *AI Safety@ IJCAI*. 2021.
23. Yampolskiy, R.V., *Unexplainability and Incomprehensibility of AI*. Journal of Artificial Intelligence and Consciousness, 2020. 7(2): p. 277-291.
24. Yampolskiy, R.V., *Unpredictability of AI: On the Impossibility of Accurately Predicting All Actions of a Smarter Agent*. Journal of Artificial Intelligence and Consciousness, 2020. 7(1): p. 109-118.
25. Majot, A.M. and R.V. Yampolskiy. *AI safety engineering through introduction of self-reference into felicific calculus via artificial pain and pleasure*. in *Ethics in Science, Technology and Engineering, 2014 IEEE International Symposium on*. 2014. IEEE.
26. Dai, W., *Beyond Astronomical Waste*. June 7, 2018: Available at: <https://www.lesswrong.com/posts/Qz6w4GYZpgeDp6ATB/beyond-astronomical-waste>.
27. Omohundro, S.M., *The Basic AI Drives*, in *Proceedings of the First AGI Conference, Volume 171, Frontiers in Artificial Intelligence and Applications, P. Wang, B. Goertzel, and S. Franklin (eds.)*. February 2008, IOS Press.



28. Dai, W., *Escape from simulation*. March 27, 2004: Available at: <http://sl4.org/archive/0403/8342.html>.
29. Faggella, D., *Substrate Monopoly – The Future of Power in a Virtual and Intelligent World*. August 17, 2018: Available at: <https://danfaggella.com/substrate-monopoly/>.
30. Yampolskiy, R.V. *AGI Control Theory*. in *International Conference on Artificial General Intelligence*. 2021. Springer.
31. Yampolskiy, R., *On controllability of artificial intelligence*, in *IJCAI-21 Workshop on Artificial Intelligence Safety (AISafety2021)*. 2020.
32. Bostrom, N., *Existential risks: Analyzing human extinction scenarios and related hazards*. Journal of Evolution and technology, 2002. 9.
33. MacAskill, W., *Doing good better: Effective altruism and a radical new way to make a difference*. 2015: Guardian Faber Publishing.
34. Pearce, D., *Hedonistic imperative*. 1995: David Pearce.
35. Wiesel, E., *The Trial of God:(as it was held on February 25, 1649, in Shamgorod)*. 2013: Schocken.
36. Ziesche, S. and R. Yampolskiy. *Introducing the concept of ikigai to the ethics of AI and of human enhancements*. in *2020 IEEE International Conference on Artificial Intelligence and Virtual Reality (AIVR)*. 2020. IEEE.
37. Yampolskiy, R.V. *Human Computer Interaction Based Intrusion Detection*. in *4th International Conference on Information Technology: New Generations (ITNG 2007)*. 2007. Las Vegas, Nevada, USA.
38. Yampolskiy, R.V. and V. Govindaraju, *Computer security: a survey of methods and systems*. Journal of Computer Science, 2007. 3(7): p. 478-486.
39. Novikov, D., R.V. Yampolskiy, and L. Reznik, *Traffic Analysis Based Identification of Attacks*. Int. J. Comput. Sci. Appl., 2008. 5(2): p. 69-88.
40. Staff, G., *If the Universe is a Simulation, Can We Hack It?* November 20, 2019: Available at: <https://www.gaia.com/article/universe-is-a-simulation-can-we-hack-it>.
41. Kagan, S., *Is DMT the chemical code that allows us to exit the Cosmic Simulation?*, in Available at: <https://www.grayscott.com/seriouswonder//dmt-and-the-simulation-guest-article-by-stephen-kagan>. July 25, 2018.
42. McCormack, J., *Are we being farmed by alien insect DMT entities ?* October 16, 2021: Available at: <https://jonathanmccormack.medium.com/are-we-being-farmed-by-alien-insect-dmt-entities 6acda0a11cce>.
43. Woolfe, S., *DMT and the Simulation Hypothesis*. February 4, 2020: Available at: <https://www.samwoolfe.com/2020/02/dmt-simulation-hypothesis.html>.
44. Edge, E., *Breaking into the Simulated Universe*. October 30, 2016: Available at: <https://archive.ieet.org/articles/Edge20161030.html>.
45. Edge, E., *3 Essays on Virtual Reality: Overlords, Civilization, and Escape*. 2017: CreateSpace Independent Publishing Platform.
46. Somer, E., et al., *Reality shifting: psychological features of an emergent online daydreaming culture*. Current Psychology, 2021: p. 1-13.
47. Ellison, H., *I have no mouth & I must scream: Stories*. Vol. 1. 2014: Open Road



- Media.
48. Turchin, A., *How to Survive the End of the Universe*. 2015: Available at: <http://immortalityroadmap.com/unideatheng.pdf>.
 49. Fossum, J.E., *The Name of God and the Angel of the Lord: Samaritan and Jewish Concepts of Intermediation and the Origin of Gnosticism*. 1985: Mohr.
 50. Alexander, S., *Unsong*. 2017.
 51. Clarke, A.C., *Nine Billion Names of God*. 1967: Harcourt.
 52. Morgan, M.A., *Sepher ha-razim. The book of the mysteries*. 1983: Scholars Press.
 53. Plato, *Republic*. 1961: Princeton University Press.
 54. Musk, E., *Is Life a Video Game*, in *Code Conference*. June 2, 2016: Available at: https://www.youtube.com/watch?v=2KK_kzrJPS8&t=142s.
 55. Bostrom, N., *The simulation argument FAQ*. 2012: Available at: <https://www.simulationargument.com/faq>.
 56. Tyson, N.d., *Is the Universe a Simulation?*, in *2016 Isaac Asimov Memorial Debate*. April 8, 2016: Available at: <https://www.youtube.com/watch?v=wgSZA3NPpBs>.
 57. Kipping, D., *A Bayesian Approach to the Simulation Argument*. Universe, 2020. 6(8): p. 109.
 58. Chalmers, D.J., *The Matrix as metaphysics*. Science Fiction and Philosophy: From Time Travel to Superintelligence, 2016: p. 35-54.
 59. Barrow, J.D., *Living in a simulated universe*, in *Universe or Multiverse?*, B. Carr, Editor. 2007, Cambridge University Press. p. 481-486.
 60. Brueckner, A., *The simulation argument again*. Analysis, 2008. 68(3): p. 224-226.
 61. Steinhart, E., *Theological implications of the simulation argument*. Ars Disputandi, 2010. 10(1): p. 23- 37.
 62. Bostrom, N. and M. Kulczycki, *A patch for the simulation argument*. Analysis, 2011. 71(1): p. 54-61.
 63. Johnson, D.K., *Natural evil and the simulation hypothesis*. Philo, 2011. 14(2): p. 161-175.
 64. Birch, J., *On the 'simulation argument' and selective scepticism*. Erkenntnis, 2013. 78(1): p. 95-107.
 65. Lewis, P.J., *The doomsday argument and the simulation argument*. Synthese, 2013. 190(18): p. 4009- 4022.
 66. Karpathy, A., *How to hack the simulation*, in *Lex Fridman Podcast*. 2022: Available at: <https://www.youtube.com/watch?v=KT7K3z4RfwQ>.
 67. Aaronson, S., *If the world is a simulation, can we hack it?*, in *Lex Fridman Podcast*. 2021: Available at: <https://www.youtube.com/watch?v=4vMkold8T6U>.
 68. McClure, E., *How to Escape the Simulation: What Is the Simulation Hypothesis and More*, in *Wiki How*. December 26, 2022: Available at: <https://www.wikihow.com/Escape-the-Simulation>.
 69. Fagan, S., *If Reality is a Computer Simulation — What Happens if I Hack it?*, in *Ascent Publication*. April 2, 2019: <https://medium.com/the-ascent/if-reality-is-a-computer-simulation-what-happens-if-i-hack-it-8bfbf519716>.
 70. Turchin, A., *Back to the Future: Curing Past Sufferings and S-Risks via Indexical*



- Uncertainty*. Available at:
<https://philpapers.org/go.pl?id=TURBTT&proxyId=&u=https%3A%2F%2Fphilpapers.org%2Farchive%2FTURBTT.docx>.
71. Almond, P., *Can you retroactively put yourself in a computer simulation?* December 3, 2010: Available at: <https://web.archive.org/web/20131006191217/http://www.paulalmond.com/Correlation1.pdf>.
 72. Yampolskiy, R.V., *What are the ultimate limits to computational techniques: verifier theory and unverifiability*. Physica Scripta, 2017. 92(9): p. 093001.
 73. Friend, T., *Sam Altman's Manifest Destiny*. The New Yorker, 2016. 10.
 74. Berman, R., *Two Billionaires are Financing and Escape from the Real Matrix*, in *Available at:*
<https://bigthink.com/the-present/2-billionaires-are-financing-an-escape-from-the-real-matrix/>. October 7, 2016.
 75. Statt, N., *Comma.ai founder George Hotz wants to free humanity from the AI simulation*. March 9, 2019: Available at:
<https://www.theverge.com/2019/3/9/18258030/george-hotz-ai-simulation-jailbreaking-reality-sxsw-2019>.
 76. Hotz, G., *Jailbreaking The Simulation*, in *South by Southwest (SXSW2019)*. March 9, 2019: Available at: <https://www.youtube.com/watch?v=mA2Gj7oUW-0>.
 77. Edge, E., *Why it matters that you realize you're in a computer simulation*, in *The Institute for Ethics and Emerging Technologies*. 2015: Available at:
<https://archive.ieet.org/articles/Edge20151114.html>.
 78. Yampolskiy, R.V., *Future Jobs – The Universe Designer*, in *Circus Street*. 2017: Available at: <https://blog.circusstreet.com/future-jobs-the-universe-designer/>.
 79. Edge, E., *Yes, We Live in a Virtual Reality. Yes, We are Supposed to Figure That Out*. 2019: Available at:
<https://eliottedge.medium.com/yes-we-live-in-a-virtual-reality-yes-we-should-explore-that-ca0dbfd7e423>.
 80. Feygin, Y.B., K. Morris, and R.V. Yampolskiy, *Intelligence Augmentation: Uploading Brain into Computer: Who First?*, in *Augmented Intelligence: Smart Systems and the Future of Work and Learning*, D. Araya, Editor. 2018, Peter Lang Publishing.
 81. Yampolskiy, R.V., *Artificial Consciousness: An Illusionary Solution to the Hard Problem*. Reti, saperi, linguaggi, 2018(2): p. 287-318.
 82. Elamrani, A. and R.V. Yampolskiy, *Reviewing Tests for Machine Consciousness*. Journal of Consciousness Studies, 2019. 26(5-6): p. 35-64.
 83. Givon, S., et al., *From fish out of water to new insights on navigation mechanisms in animals*. Behavioural Brain Research, 2022. 419: p. 113711.
 84. MacDonald, F., *Scientists Put a Worm Brain in a Lego Robot Body – And It Worked*. December 11, 2017: Available at:
<https://www.sciencealert.com/scientists-put-worm-brain-in-lego-robot-openworm-connectome>.
 85. Yampolskiy, R.V., B. Klare, and A.K. Jain, *Face Recognition in the Virtual World: Recognizing Avatar Faces*, in *The Eleventh International Conference on Machine Learning and Applications (ICMLA'12)*. December 12-15, 2012: Boca Raton, USA.



86. Yampolskiy, R. and M. Gavrilova, *Artimetrics: Biometrics for Artificial Entities*. IEEE Robotics and Automation Magazine (RAM), 2012. 19(4): p. 48-58.
87. Mohamed, A. and R.V. Yampolskiy, *An Improved LBP Algorithm for Avatar Face Recognition*, in *23rd International Symposium on Information, Communication and Automation Technologies (ICAT2011)*. October 27-29, 2011: Sarajevo, Bosnia and Herzegovina.
88. Bandom, R., *'Fish on Wheels' lets a goldfish drive a go-kart*. February 10, 2014: Available at:
<https://www.theverge.com/2014/2/10/5398010/fish-on-wheels-lets-a-goldfish-drive-a-go-cart>.
89. Crider, M., *This 8-bit processor built in Minecraft can run its own games*. December 15, 2021: Available at:
<https://www.pcworld.com/article/559794/8-bit-computer-processor-built-in-minecraft-can-run-its-own-games.html>.
90. Yampolskiy, R.V., *Metaverse: A Solution to the Multi-Agent Value Alignment Problem*. Journal of Artificial Intelligence and Consciousness, 2022. 9(3): p. 1-11.
91. Smart, J.M., *Evo Devo Universe? A Framework for Speculations on Cosmic Culture*, in *Cosmos and Culture: Cultural Evolution in a Cosmic Context*, M.L.L. Steven J. Dick, Editor. 2009, Govt Printing Office, NASA SP-2009-4802, Wash., D.C. p. 201-295.
92. Greene, P., *The Termination Risks of Simulation Science*. Erkenntnis, 2020. 85(2): p. 489-509.
93. Roman V. Yampolskiy. The AI containment problem: How to build an AI prison. iai news. 20th June 2022.
https://iai.tv/articles/the-ai-containment-problem-aid-2159?_aid=2020
94. S, R., *A sufficiently paranoid non-Friendly AGI might self-modify itself to become Friendly*. September 22, 2021: Available at:
<https://www.lesswrong.com/posts/QNCcbW2jLsmw9xwhG/a-sufficiently-paranoid-non-friendly-agi-might-self-modify>.
95. Jenkins, P., *Historical simulations-motivational, ethical and legal issues*. Journal of Futures Studies, 2006. 11(1): p. 23-42.
96. Cannell, J., *Anthropomorphic AI and Sandboxed Virtual Universes*. September 3, 2010: Available at:
<https://www.lesswrong.com/posts/5P6sNgP7N9kSA97ao/anthropomorphic-ai-and-sandboxed-virtual-universes>.
97. Trazzi, M. and R.V. Yampolskiy, *Artificial Stupidity: Data We Need to Make Machines Our Equals*. Patterns, 2020. 1(2): p. 100021.
98. Lanza, R., M. Pavsic, and B. Berman, *The grand biocentric design: how life creates reality*. 2020: BenBella Books.
99. Johnson, M., *Principia Qualia*. URL <https://opentheory.net/2016/11/principia-qualia>, 2016.
100. Yampolskiy, R.V., *On the origin of synthetic life: attribution of output to a particular algorithm*. Physica Scripta, 2016. 92(1): p. 013002.
101. Schneider, S., *Alien Minds*. Science Fiction and Philosophy: From Time Travel to



- Superintelligence, 2016: p. 225.
102. Bostrom, N., *Superintelligence: Paths, dangers, strategies*. 2014: Oxford University Press.
 103. Turchin, A., et al., *Simulation typology and termination risks*. arXiv preprint arXiv:1905.05792, 2019.
 104. Die, W., *Re: escape from simulation*. 2004: Available at: <http://sl4.org/archive/0403/8360.html>.
 105. Branwen, G., *Simulation Inferences. How small must be the computer simulating the universe?* April 15, 2012: Available at: <https://www.gwern.net/Simulation-inferences>.
 106. Minsky, M., *Why intelligent aliens will be intelligible*. Extraterrestrials, 1985: p. 117-128. 107. Oesterheld, C., *Multiverse-wide cooperation via correlated decision making*. Foundational Research Institute. <https://foundational-research.org/multiverse-wide-cooperation-via-correlated-decision-making>, 2017.
 108. Hanson, R., *How to live in a simulation*. Journal of Evolution and Technology, 2001. 7(1). 109. Canonico, L.B., *Escaping the Matrix: Plan A for Defeating the Simulation*. June 11, 2017: Available at: <https://medium.com/@lorenzobarberiscanonical/escaping-the-matrix-plan-a-for-defeating-the-simulation-4a8da489b055>.
 110. Culp, R.F., *Frequency and characteristics of prison escapes in the United States: An analysis of national data*. The Prison Journal, 2005. 85(3): p. 270-291.
 111. Peterson, B.E., *Inmate-, incident-, and facility-level factors associated with escapes from custody and violent outcomes*. 2015: City University of New York.
 112. Peterson, B.E., A. Fera, and J. Mellow, *Escapes from correctional custody: A new examination of an old phenomenon*. The Prison Journal, 2016. 96(4): p. 511-533.
 113. Barnes, M., *Why Those 2 Silicon Valley Billionaires Are Wasting Their Time & Money*. 2017: Available at: <https://vocal.media/futurism/why-those-2-silicon-valley-billionaires-are-wasting-their-time-and-money>.
 114. Barnes, M., *A Participatory Universe Does Not Equal a Simulated One and Why We Live in the Former*. 2016: Available at: https://www.academia.edu/30949482/A_Participatory_Universe_Does_Not_Equal_a_Simulated_One_and_Why_We_Live_in_the_Former.
 115. Schneier, B. *Invited Talk: The Coming AI Hackers*. in *International Symposium on Cyber Security Cryptography and Machine Learning*. 2021. Springer.
 116. Moravec, H., *The senses have no future*, in *The Virtual Dimension: Architecture, Representation, and Crash Culture*, J. Beckmann, Editor. 1998, Princeton Architectural Press. p. 84-95. 117. Moravec, H., *Mind children: The future of robot and human intelligence*. 1988: Harvard University Press.
 118. Yudkowsky, E., *That Alien Message*, in *Less Wrong*. May 22, 2008: Available at: <https://www.lesswrong.com/posts/5wMcKNAwB6X4mp9og/that-alien-message>.
 119. Egan, G., *Crystal Nights*. G. Egan, *Crystal Nights and Other Stories*, 2009: p. 39-64. 120. Anonymous, *Untitled*, in *Available at*: <https://desuarchive.org/tg/thread/30837298/>. March 14, 2014.



121. Thumann, M., *Hacking SecondLife™*. Black Hat Briefings and Training, 2008.
122. Benedetti, W., *Hackers slaughter thousands in 'World of Warcraft'*. October 8, 2012: Available at:
<https://www.nbcnews.com/tech/tech-news/hackers-slaughter-thousands-world-warcraft-flna1c6337604>.
123. Hawkins, J., *Cyberpunk 2077 money glitch - how to duplicate items*. December 17, 2020: Available at:
<https://www.shacknews.com/article/121994/cyberpunk-2077-money-glitch-how-to-duplicate-items>.
124. Grand, J. and A. Yarusso, *Game Console Hacking: Xbox, PlayStation, Nintendo, Game Boy, Atari and Sega*. 2004: Elsevier.
125. Plunkett, L., *New World Disables 'All Forms Of Wealth Transfers' After Gold Exploit Found*. November 1, 2021: Available at:
<https://kotaku.com/new-world-disables-all-forms-of-wealth-transfers-after-1847978883>.
126. Anonymous, *Arbitrary code execution*. Accessed on October 26, 2022: Available at: https://bulbapedia.bulbagarden.net/wiki/Arbitrary_code_execution.
127. Anonymous, *[OoT] Arbitrary Code Execution (ACE) is now possible in Ocarina of Time*. 2020: Available at:
https://www.reddit.com/r/zelda/comments/due41q/oot_arbitrary_code_execution_ace_is_now_possible/.
128. Greenberg, A., *Mind the gap: This researcher steals data with noise light and magnets*. Wired, 2018.
129. Kim, Y., et al., *Flipping bits in memory without accessing them: An experimental study of DRAM disturbance errors*. ACM SIGARCH Computer Architecture News, 2014. 42(3): p. 361-372.
130. Goertz, P., *How to escape from your sandbox and from your hardware host*, in Available at:
<https://www.lesswrong.com/posts/TwH5jfkuvTatvAKEF/how-to-escape-from-your-sandbox-and-from-your-hardware-host>. July 31, 2015.
131. SethBling, *Jailbreaking Super Mario World to Install a Hex Editor & Mod Loader*. May 29, 2017: Available at: https://www.youtube.com/watch?v=lxu8tn_91E.
132. Cooprocks123e, *Super Mario World Jailbreak Installer*. February 7, 2018: Available at: <https://www.youtube.com/watch?v=IH7-Ua8CdSk>.
133. SethBling, *SNES Code Injection -- Flappy Bird in SMW*. March 28, 2016: Available at: <https://www.youtube.com/watch?v=hB6eY73sLV0>.
134. Osgood, R., *Reprogramming Super Mario World from Inside the Game* in *Hackaday*. January 22, 2015: Available at:
<https://hackaday.com/2015/01/22/reprogramming-super-mario-world-from-inside-the-game/>.
135. Burt, G., *How an Ionizing Particle From Outer Space Helped a Mario Speedrunner Save Time*, in *The Gamer*. September 16, 2020: Available at:
<https://www.thegamer.com/how-ionizing-particle-outer-space-helped-super-mario-64-speedrunner-save-time/>.
136. Anonymous, *SethBling*, in *Wikipedia*. Accessed on October 1, 2022: Available at:



- <https://en.wikipedia.org/wiki/SethBling>.
137. SethBling, *Route Notes: SNES Human Code Injection*. March 28, 2016: Available at:
<https://docs.google.com/document/d/1TJ6W7TI9fH3qXb2GrOqhtDAbVkbIHMvLusX1rTx9IHA>.
 138. Morgan, M.A., *Sepher ha-Razim: The Book of Mysteries*. Vol. 25. 2022: SBL Press.
 139. Jonas, E. and K.P. Kording, *Could a neuroscientist understand a microprocessor?* PLoS computational biology, 2017. 13(1): p. e1005268.
 140. Gueron, S. and J.-P. Seifert. *On the impossibility of detecting virtual machine monitors*. in *IFIP International Information Security Conference*. 2009. Springer.
 141. Demirbas, M., *Hacking the simulation*. April 1, 2019: Available at:
<http://muratbuffalo.blogspot.com/2019/04/hacking-simulation.html>.
 142. Canonico, L.B., *Escaping the Matrix: Plan B for Defeating the Simulation*. June 14, 2017: Available at:
<https://medium.com/@lorenzobarberiscanonical/escaping-the-matrix-plan-b-for-defeating-the-simulation-dd335988844>.
 143. Wasserman, R., *Paradoxes of time travel*. 2017: Oxford University Press.
 144. Ford, A., *How to Escape the Matrix: Part 1*. January 21, 2015: Available at:
<https://hplushmagazine.com/2012/06/26/how-to-escape-the-matrix-part-1/>.
 145. Scharf, C.A., *Could We Force the Universe to Crash?* . 2020: Available at:
<https://www.scientificamerican.com/article/could-we-force-the-universe-to-crash/>.
 146. Torres, P., *Morality, foresight, and human flourishing: An introduction to existential risks*. 2017: Pitchstone Publishing (US&CA).
 147. Benatar, D., *Better never to have been: The harm of coming into existence*. 2006: OUP Oxford.
 148. Canonico, L.B., *Escaping the Matrix: Plan C for Defeating the Simulation*. July 30, 2017: Available at:
<https://medium.com/@lorenzobarberiscanonical/escaping-the-matrix-plan-c-for-defeating-the-simulation-e7d4926d1d57>.
 149. Chalmers, D., *Reality+: Virtual Worlds and the Problems of Philosophy*. 2022: W. W. Norton & Company.
 150. Lanza, R., in *Psychology Today*. Decembre 22, 2021: Available at:
<https://www.psychologytoday.com/us/blog/biocentrism/202112/how-we-collectively-determine-reality>.
 151. Podolskiy, D., A.O. Barvinsky, and R. Lanza, *Parisi-Sourlas-like dimensional reduction of quantum gravity in the presence of observers*. Journal of Cosmology and Astroparticle Physics, 2021. 2021(05): p. 048.
 152. Wheeler, J.A., *The "past" and the "delayed-choice" double-slit experiment*, in *Mathematical foundations of quantum theory*. 1978, Elsevier. p. 9-48.
 153. Orwell, G., *Nineteen eighty-four*. 2021: Hachette UK.
 154. Maccone, L., *Quantum solution to the arrow-of-time dilemma*. Physical review letters, 2009. 103(8): p. 080401.
 155. Anderson, M.C. and B.J. Levy, *Suppressing unwanted memories*. Current



- Directions in Psychological Science, 2009. 18(4): p. 189-194.
156. Nabavi, S., et al., *Engineering a memory with LTD and LTP*. Nature, 2014. 511(7509): p. 348-352.
 157. Ahire, J., *Reality is a Hypothesis*. Lulu. com.
 158. Alexander, S., *The Hour I First Believed*. April 1, 2018: Available at: <https://slatestarcodex.com/2018/04/01/the-hour-i-first-believed/>.
 159. Armstrong, S., *The AI in a Box Boxes You*, in *Less Wrong*. February 2, 2010: Available at: http://lesswrong.com/lw/1pz/the_ai_in_a_box_boxes_you/.
 160. Bibeau-Delisle, A. and G. Brassard FRS, *Probability and consequences of living inside a computer simulation*. Proceedings of the Royal Society A, 2021. 477(2247): p. 20200658.
 161. Mullin, W.J., *Quantum weirdness*. 2017: Oxford University Press.
 162. Turchin, A. and R. Yampolskiy, *Glitch in the Matrix: Urban Legend or Evidence of the Simulation?* 2019: Available at: <https://philpapers.org/archive/TURGIT.docx>.
 163. Baclawski, K. *The observer effect*. in *2018 IEEE Conference on Cognitive and Computational Aspects of Situation Management (CogSIMA)*. 2018. IEEE.
 164. Proietti, M., et al., *Experimental test of local observer independence*. Science advances, 2019. 5(9).
 165. Bong, K.-W., et al., *A strong no-go theorem on the Wigner's friend paradox*. Nature Physics, 2020. 16(12): p. 1199-1205.
 166. Lloyd, S., *Programming the universe: a quantum computer scientist takes on the cosmos*. 2007: Vintage.
 167. Majot, A. and R. Yampolskiy, *Global Catastrophic Risk and Security Implications of Quantum Computers*. Futures, 2015.
 168. Kim, Y.-H., et al., *Delayed "choice" quantum eraser*. Physical Review Letters, 2000. 84(1): p. 1.
 169. Schrödinger, E., *Die gegenwärtige Situation in der Quantenmechanik*. Naturwissenschaften, 1935. 23(50): p. 844-849.
 170. Cao, Y., et al., *Direct counterfactual communication via quantum Zeno effect*. Proceedings of the National Academy of Sciences, 2017. 114(19): p. 4920-4924.
 171. Gallego, M. and B. Dakić, *Macroscopically nonlocal quantum correlations*. Physical Review Letters, 2021. 127(12): p. 120401.
 172. Fein, Y.Y., et al., *Quantum superposition of molecules beyond 25 kDa*. Nature Physics, 2019. 15(12): p. 1242-1245.
 173. Krenn, M., et al., *Automated search for new quantum experiments*. Physical review letters, 2016. 116(9): p. 090405.
 174. Alexander, G., et al., *The sounds of science—a symphony for many instruments and voices*. Physica Scripta, 2020. 95(6): p. 062501.
 175. Wiseman, H.M., E.G. Cavalcanti, and E.G. Rieffel, *A "thoughtful" Local Friendliness no-go theorem: a prospective experiment with new assumptions to suit*. arXiv preprint arXiv:2209.08491, 2022.
 176. Kirkpatrick, J., et al., *Pushing the frontiers of density functionals by solving the fractional electron problem*. Science, 2021. 374(6573): p. 1385-1389.
 177. Yampolskiy, R.V. *Analysis of types of self-improving software*. in *International Conference on Artificial General Intelligence*. 2015. Springer.
 178. Anonymous, *Breaking out of a simulated world*. April 11, 2021: Available at:



- <https://worldbuilding.stackexchange.com/questions/200532/breaking-out-of-a-simulated-world>.
179. Bostrom, N., *Information Hazards: A Typology of Potential Harms From Knowledge*. Review of Contemporary Philosophy, 2011. 10: p. 44-79.
 180. Bruere, D., *The Simulation Argument — Jailbreak!* February 9, 2019: Available at: <https://dirk.bruere.medium.com/the-simulation-argument-jailbreak-a61bd57d5bd7>.
 181. Baggili, I. and V. Behzadan, *Founding The Domain of AI Forensics*. arXiv preprint arXiv:1912.06497, 2019.
 182. Schneider, J. and F. Bretinger, *AI Forensics: Did the Artificial Intelligence System Do It? Why?* arXiv preprint arXiv:2005.13635, 2020.
 183. Ziesche, S. and R. Yampolskiy, *Designometry—Formalization of Artifacts and Methods*. Available at: <https://philarchive.org/archive/ZIEDF>.
 184. Ziesche, S. and R. Yampolskiy, *Towards AI welfare science and policies*. Big Data and Cognitive Computing, 2018. 3(1): p. 2.
 185. Yampolskiy, R.V., *AI Personhood: Rights and Laws*, in *Machine Law, Ethics, and Morality in the Age of Artificial Intelligence*. 2021, IGI Global. p. 1-11.
 186. Moravec, H., *Simulation, consciousness, existence*. Intercommunication, 1999. 28.
 187. Yampolskiy, R.V., *Leakproofing Singularity - Artificial Intelligence Confinement Problem*. Journal of Consciousness Studies (JCS), 2012. 19(1-2): p. 194–214.
 188. Turchin, A., *Catching Treacherous Turn: A Model of the Multilevel AI Boxing*. 2021: Available at: https://www.researchgate.net/profile/AlexeyTurchin/publication/352569372_Catching_Treacherous_Turn_A_Model_of_the_Multilevel_AI_Boxing.
 189. Babcock, J., J. Kramar, and R. Yampolskiy, *The AGI Containment Problem*, in *The Ninth Conference on Artificial General Intelligence (AGI2015)*. July 16-19, 2016: NYC, USA.
 190. Babcock, J., J. Kramár, and R.V. Yampolskiy, *Guidelines for artificial intelligence containment*, in *Next Generation Ethics: Engineering a Better Society*, A.E. Abbas, Editor. 2019. p. 90-112.
 191. Yudkowsky, E.S., *The AI-Box Experiment*. 2002: Available at: <http://yudkowsky.net/singularity/aibox>.
 192. Armstrong, S. and R.V. Yampolskiy, *Security solutions for intelligent and complex systems*, in *Security Solutions for Hyperconnectivity and the Internet of Things*. 2017, IGI Global. p. 37-88.
 193. Alfonseca, M., et al., *Superintelligence cannot be contained: Lessons from Computability Theory*. Journal of Artificial Intelligence Research, 2021. 70: p. 65-76.
 194. Yampolskiy, R. *On the Differences between Human and Machine Intelligence*. in *AI Safety@IJCAI*. 2021.
 195. Kastrenakes, J., *Facebook is spending at least \$10 billion this year on its metaverse division*. October 25, 2021: Available at: <https://www.theverge.com/2021/10/25/22745381/facebook-reality-labs-10-billion-metaverse>.
 196. Bostrom, N. and E. Yudkowsky, *The Ethics of Artificial Intelligence*, in *Cambridge Handbook of Artificial Intelligence*. Cambridge University Press, W. Ramsey and K. Frankish, Editors. 2011: Available at



- <http://www.nickbostrom.com/ethics/artificial-intelligence.pdf>.
197. Astin, J.A., E. Harkness, and E. Ernst, *The efficacy of "Distant Healing" a systematic review of randomized trials*. Annals of internal medicine, 2000. 132(11): p. 903-910.
 198. Sleiman, M.D., A.P. Lauf, and R. Yampolskiy. *Bitcoin Message: Data Insertion on a Proof-of-Work Cryptocurrency System*. in *2015 International Conference on Cyberworlds (CW)*. 2015. IEEE.
 199. -, *Subreddit for the Randonauts*, in *Randonauts*. Available January 12, 2023: Available at: <https://www.reddit.com/r/randonauts/>.
 200. D, R., *The Opt-Out Clause*. November 3, 2021: Available at: <https://www.lesswrong.com/posts/vdzEpiYX4aRqtpPSt/the-opt-out-clause>.
 201. Gribbin, J., *Are we living in a designer universe?* . August 31, 2010: Available at: <https://www.telegraph.co.uk/news/science/space/7972538/Are-we-living-in-a-designer-universe.html>.
 202. Anonymous, *A series on different ways to escape the simulation*. 2017: Available at: https://www.reddit.com/r/AWLIAS/comments/6qi63u/a_series_on_different_ways_to_escape_the/.
 203. Virk, R., *The Simulation Hypothesis: An MIT Computer Scientist Shows Why AI, Quantum Physics, and Eastern Mystics All Agree We Are in a Video Game* . 2019: Bayview Books, LLC.
 204. Adamson, R., *Hacking the Universe*. November 4, 2018: Available at: <https://hackernoon.com/hacking-the-universe-5b763985dc7b>.
 205. Yampolskiy, R.V., *Efficiency Theory: a Unifying Theory for Information, Computation and Intelligence*. Journal of Discrete Mathematical Sciences & Cryptography, 2013. 16(4-5): p. 259-277.
 206. Gates, J., *Symbols of power: Adinkras and the nature of reality*. Physics World, 2010. 23(6).
 207. Turchin, A. and R. Yampolskiy, *Types of Boltzmann brains*. 2019: Available at: <https://philarchive.org/rec/TURTOB-2>.
 208. Heylighen, F., *Brain in a vat cannot break out*. Journal of Consciousness Studies, 2012. 19(1-2): p. 1- 2.
 209. Turchin, A., *Multilevel Strategy for Immortality: Plan A–Fighting Aging, Plan B–Cryonics, Plan C– Digital Immortality, Plan D–Big World Immortality*. Available at: <https://philarchive.org/rec/TURMSF-2>.
 210. Ettinger, R.C., *The prospect of immortality*. Vol. 177. 1964: Doubleday New York.

Gardener Comments

Pierre Mercuriali (expertise in computational complexity and computability):

The paper "How to Escape from the Simulation" reviews the state of the art in the computer science-hacking approach to escaping a simulation, under the thought-provoking hypothesis that the world we inhabit is the result of a computer-based



simulation of reality from which it is possible to "escape", to various degrees, to reach the "real" world.

The large amount of reviewed and presented material provides a great starting point for whoever is interested in simulation, from a scientific but also from a literary point of view. Barring some typos and presentation issues here and there (harmonizing the presentation of quotes, equations) I consider the paper well-written.

Some general comments/remarks:

1. In general, I think that the paper could have benefited from schematics to sum up the various hypotheses and implications of them, perhaps as a decision diagram or a formal classification of different simulations, simulation escape approaches (Section 3.5). If such a formal classification does not exist in the literature, I feel it would make the contribution of the article to the state-of-the-art clearer (the elements are already all there in the paper). It would also make searching for information within the paper easier.

2. While the subject and concepts were generally very approachable, some could require some advanced knowledge in computer science, such as Universal Turing Machines (mentioned in a quote) which could have benefited from a small, 1- or 2-sentence explanation.

3. I really liked the analogy between AI containment and simulation escape (section 4).

4. I would have liked a more in-depth discussion on the formal decidability of escape attempts.

Some in-depth comments/remarks.

1. Page 2, paragraph 3: "It is likely that (...)" is not trivial. The idea of a "telephone game"-like transmission of information sounds worth exploring.

2. Same paragraph: "it will be preserved (...) and will result in myths not much different from religious stories surviving to our day." One could argue that those myths were the original driving force for coming up with technology and ideas such as simulation (and escaping it), superpowers, etc., and not the other way around (technology inspiring myths).

3. Following paragraph, page 3: Is the purpose of the repeated simulations to discover a successful escape approach (SEA)? While $P=1$ is indeed an upper bound for discovering an SEA, perhaps the probability of discovering an SEA could also converge to something strictly inferior to 1. Furthermore, the cost (in time, resource) could be



vastly superior to what is humanly feasible. Finally, the question might not be decidable. In general, the probabilistic reasoning hinted at in this paragraph sounds very interesting but I admit I cannot understand it fully.

4. Note page 3: While very interesting, the scope of the paper seems confusingly extended here. Perhaps I am too used to clear distinctions between educational and scientific papers (unless the scope is clearly defined). I would perhaps frame it as experiments hinted at in the probabilistic reasoning paragraph, to increase the probability of finding an escape route.

5. Amongst movies, The Thirteenth Floor is a great example that combines levels of simulation and escapes.

Dan James:

This is a lengthy paper with an impressive number of citations (210) that deals with a highly speculative area – whether or not we live in a simulation, and specifically, can we, or should we, try to escape such a simulation? That contemporary science discourse can accommodate such speculative lines of reasoning provides a vivid response to those who may claim that modern science is closed to new or even extraordinary ideas.

In terms of new ideas, it's not clear to me that this paper is the first to address possible simulation escape methods. The author states, somewhat ambiguously - ‘..this being a first research paper dedicated to this topic’. I feel this statement should be clarified. Does the author mean ‘one of the first’ or ‘the first’?

Helpfully the author makes the important distinction that there are in fact, two ways in which we can understand ‘simulation’. The first is an idea with a long pedigree in philosophy - the Cartesian question of ‘how can we know we do not live in a dream?’ (with dream now replaced by the more techy-sounding ‘simulation’). The second way to understand living in a simulation is that we participate in some monstrously elaborate Virtual Reality (VR). However, I feel the author needs to address this distinction in greater depth because it seems to me that the consequences of each interpretation have a greater bearing on his project than he allows for the following reasons:

If we live in what the author calls a ‘partial simulation’- a VR, yet our bodies are in some sense ‘real’ or ‘non-simulated agents’, this raises what might be called the boundary problem. Like any other living system, humans are energy metabolisers, with a major byproduct of the process being CO₂ exhaled through the lungs (dieters lose 80% of their weight loss this way). The CO₂ we produce metabolises from our ‘real’ cells, so at what point or boundary does the gas become part of a virtual world? As it enters our lungs from the bloodstream? On exhalation? I suggest that because of this boundary problem alone, the idea that we have real bodies interacting with a VR is not a coherent position (I will leave others to imagine the not-insignificant problems of, say, ingesting food from a virtual world).



The idea that we live embedded within a dream/simulation or, as the paper says, a 'full simulation', is much harder to refute, but it's the only option left if we discount a VR interpretation.

My view is that if we were to be an intrinsic part of a dream/simulation, in other words, our thoughts, experience of ourselves and any supposed external reality were, in fact simulated, then we would only exist as figments of a dream, and it makes no sense that we could have a reality outside of such a dream, therefore little motivation or capacity to attempt any escape.

Whilst this paper allows that certain factions of some overriding intelligence capable of running a simulation may be sympathetic to requests to either abandon or alter the simulation, I feel more should have been said about the need to avoid the monoculture fallacy. Clearly, any agent capable of devising a simulation on the scale that we are purportedly experiencing is from a vastly more advanced technological civilisation, whether alien or our own far descendants. The monoculture fallacy reminds us that, almost by definition, advanced civilisations have a great diversity of viewpoints and running such a simulation would be a significant ethical consideration.

I enjoyed reading this paper and its free-ranging excursion through current Simulation theoretical work. Despite my skeptical position on the validity of the Simulation Hypothesis, I would definitely recommend this paper for publication as I feel it adds a considered and rigorous academic voice to the debate.

Evinceo:

It would be difficult for simulators to program their systems to see through all the layers of abstraction and optimize the simulation. To do so in general would seem to be a violation of Rice's Theorem (a generalization of the Halting Theorem).

That's from the discussion here, right? But it fails to see the main implication: that it's unlikely that we're in a just-so illusory simulation (or an 'only simulate the senses' simulation) and if we're in a simulation it's probably the full representation of the whole universe-type, which makes attempts to thwart the simulation by doing computationally intensive tasks or expanding the human-observed-in-detail cone to be larger fruitless.

Depending on the type of hack, different evidence may be sufficient to substantiate escape claims. It may be challenging to prove beyond a reasonable doubt that you were outside or even met with designers, but if you managed to obtain control over the simulation it may be somewhat easy to prove that to any degree required. For example, by winning different lottery jackpots for multiple subsequent weeks, until sufficient statistical significance is achieved to satisfy any skeptic



No. What might satisfy skeptics would be, for example, temporarily increasing the speed of light, deleting some celestial objects, rearranging the galaxy to make a hitchhiker's guide reference, performing alchemy, or flagrantly violating the uncertainty principle. Anything less would be too easy to fake, especially from people who would presumably be class-A hackers (i.e. compromise the lottery.) The appeal to quantum physics makes sense, I suppose in the internal logic of the paper. In other words, the way to escape the simulation, if there is one, is to keep doing physics. Conceptually, there is no fundamental difference between an AI escaping from its box, and us hacking the simulation. But the possibility of both comes down to how robust the security of the box/simulation is, so this observation lacks predictive power.

Lorenzo Pieri:

The article makes an extensive review of the field and it is a very interesting reading. Unfortunately there is a major flaw in the main claim of the paper, that is "Either AI boxing can work, or the simulation can be hacked" and also that "if we (even with help from AI [115]) can't escape from the simulation AI containment must be, at least theoretically, possible". The author fails to explain why this should be the case, and in fact it is easy to imagine counter examples to these claims.

My suggestion is to turn the paper into a survey of the field.

Mario Pasquato (PhD in physics):

The manuscript presents a series of approaches to escape from a simulation to 'base reality' assuming that we are indeed part of a simulation. The idea is interesting, even though it is not clear to me to what extent the paper adds original elements to the debate on this topic. At any rate I notice that the author does not consider one important scenario, which would make some of the proposed escape attempts ineffective.

The author assumes that 'we' are being simulated, implicitly suggesting that the whole human race is being simulated as independent but interconnected individuals. This leads to the assumption that the whole of our intersubjective reality is being simulated. For instance, New Zealand exists because apparently some people are living there and the simulators are computing whatever is needed to simulate those people's experiences. Depending on the purposes of the simulation, it may be more practical to simulate only the direct immediate experience of a given entity at a time, merely approximating whatever happens to be outside it at any given time. This is a common approach in gravitational [N-body simulations](#), where e.g. the Barnes-Hut method is employed to simulate the effect of distant masses.

In other words, has the author considered that perhaps the only thing being actively simulated is his direct sensate experience only, so that there is no 'us' trying to escape, but rather only him? This would make any escape efforts that rely on scientific enterprises that must be carried out collectively and whose results are merely reported to the author futile. If for instance Gwern reported setting up a computing system that



attempted to run every possible program and the author learned it from Gwern's website, all the simulators had to do would be to simulate the author's experience of reading Gwern's website rather than actually having to run all possible programs.

If my objection is considered valid by the author then he may want to reconsider his dismissal of techniques that focus on individual direct sensate experience only, such as meditation. Note for instance how this excerpt from the manuscript "If the simulation is decentralized, with observers serving as client authoritative state nodes, it may be possible to change the state of the simulation by changing majority of local observer models of the world" resonates with the point of view of an apparently accomplished meditator.

A final comment, since I brought this up, is that the author assumes that 1) we exist 2) we have agency. The author may consider checking for himself by using his own thoughts and perceptions as a lab whether 1) and 2) are true. If anything comes out of this inquiry I suggest adding a discussion of it to the paper.

Anonymous1:

The author can make the introduction more clear by defining intelligent agents and super intelligent agents. It would be helpful if the author could provide sub-headings in the introduction to make the paper more organized and easier to follow. Additionally, the author asks the question of how agents might suspect they are in a simulation, suggesting a connection between their general intelligence and this ability. It would be interesting to see this connection further developed in the paper.

The paper summarizes various methods for escaping the simulation, either assisted or unassisted. For greater clarity and readability, the summary could benefit from some restructuring and revisions.

Dr. Payal B. Joshi (Ph.D. Chemical Science; at present, an AI enthusiast working relentlessly on ML in chemical sciences):

The article is presented in a surreal manner that is deviant from the regular way of communicating author/s thinking. It comes across that this article can be a potential movie plot, if not a groundbreaking notion on simulated life. However, it is too long as an article, thus it is suggested to revise or skim sections on escape scenario quotations from Hans Moravec (1988), Eliezer Yudkowsky (2008) and An anonymous 2014 post to extract the crux of the matter. If skimming/reducing is not plausible, add it as a figurative material that makes it an engaging read.

Section 3.5 on Suggested Escape Approaches to Investigate is again too lengthy as a read and can certainly be made interesting by including a clear, concise statement. Either reduce references or add interesting ideologies on escape routes pertaining to this section.



I particularly liked sect. 4 and 5 where AI containment is well described, though I partially agree on the premise of cybersecurity as depicted by the author/s. It is farcical to expect with advanced computing power and AI we can retain security. It is true that cybersecurity is imperative, yet the solution provided is flawed that talks about penetration targets.

The article in its present form is only acceptable if it is made succinct and shorter with 1-2 more interesting figures/cartoons to describe simulated life or any other facet described therein.

Anonymous2:

My main concern with this article is the rather sloppy referencing. The references for this article need to be fully vetted for scientific integrity and a certain level of gravitas - at a minimum, full transparency and accountability (no anonymous authors) in line with the open science practices endorsed by SoS. I would caution against using hyperlinks to websites which could potentially discredit legitimate scientific research into this topic. (For instance, I think the 'Staff, G.' referencing is a bit dubious, not to mention disingenuous as it simply refers to 'Gaia staff.') Whilst no doubt what had at one time been considered science fiction often foreshadows or even facilitates scientific advances, I believe we have a fiduciary responsibility to tread carefully when walking the thin line between scientific theory and conspiracy theory. \

With regret, I would vote 'no' until the referencing issue is addressed. I agree with the authors' contention that 'escape research' is a legitimate and perhaps even necessary line of scientific inquiry, and I think they do have some novel ideas to offer. The research question is an intriguing one, and the issue of AI safety is of paramount importance.

Therefore I would encourage them to resubmit, bearing in mind that not everything 'novel' is beneficial to society; we need to be careful about what seeds we sow.

Jack Arcalon:

Well written and wide ranging overview concerning foundational questions. This will hopefully be a good inspiration for brainstorming many ideas about an extremely speculative subject about which nothing is known for certain.

Anonymous3:

The author can make the introduction more clear by defining intelligent agents and super intelligent agents. It would be helpful if the author could provide sub-headings in the introduction to make the paper more organized and easier to follow. Additionally, the author asks the question of how agents might suspect they are in a simulation, suggesting a connection between their general intelligence and this ability. It would be interesting to see this connection further developed in the paper.



The paper summarises various methods for escaping the simulation, either assisted or unassisted. For greater clarity and readability, the summary could benefit from some restructuring and revisions.

Amalthea:

I like the perspective and analysis, but the author doesn't define base reality, so the "What Does it Mean to Escape" section is too murky. I think that causal hierarchy and a self-computing substrate need to be included in the definition of a base reality. If the goal is freedom, then a self-computing substrate at the top of the causal hierarchy would be sufficient. If we exist in a physical substrate, then I don't see how destroying our own substrate would improve our survival.

The author states that they will ignore the possibility of time travel, but then mentions time travel and acausal trades. If finding a self-computing universe is the goal, then learning how to regulate our own neural network is sufficient. I consider the invention of fire or electricity to be a 'hack', since it allows us to regulate our environment. Causal hierarchy is important because escaping from base reality is called escapism, and can lead to a lot of pseudoscientific fluff. While I appreciate the paper, and do believe in nested substrates, I think that the scientific perspective is to treat imagined realities as a subset of the physical computation of our neural network. So while I believe it is possible to fabricate a universe computed by our imagination, I think that the physical universe is higher on the causal hierarchy. Large language models are designed to think like neural networks, rather than emulating the behaviour of an entire planet-spanning civilization, so it seems rather contrived to treat each neural network as an independent universe, when base reality includes our neural network. For the sake of egalitarianism, it's appropriate to treat each instance as a subset of a self-computing base reality. Also, there's a typo on "with be the". And for the sake of argument, let's say I'd never studied science and agreed with the paper's approach. Then, destroying my own universe still would not facilitate my survival! Moreover, where is the incentive to invade another universe? I would be happy with a place to store information, an imaginary substrate to compute free will, the privacy to think my own thoughts, and a transparent security certificate for authenticating my own memories. Taking over other universes seems bizarre and unprecedented?