



Simulation Expectation

Teruji Thomas¹

Received: 29 August 2023 / Accepted: 3 November 2024
© The Author(s) 2024

Abstract

I present a new argument that we are much more likely to be living in a computer simulation than in the ground-level of reality. (Similar arguments can be marshalled for the view that we are more likely to be Boltzmann brains than ordinary people, but I focus on the case of simulations.) I explain how this argument overcomes some objections to Bostrom’s classic argument for the same conclusion. I also consider to what extent the argument depends upon an internalist conception of evidence, and I refute the common line of thought that finding many simulations being run—or running them ourselves—must increase the odds that we are in a simulation.

1 Introduction

Here’s a way the world might be. At some point there exist conscious beings whose experience of the world is much like ours—let’s just call them *people*. And at some point in their history, these people run computer simulations of whole worlds, so powerful that these worlds are inhabited by other conscious beings—let’s call them *simulant people*. These simulant people again have mental lives as rich and detailed as our own. they might run further simulations on their (simulant) computers, containing other simulant people, and so on. Only if we live in the non-simulant, ground-level of reality (if there even is such a thing!) are we ourselves non-simulant people.

I will present an argument for

SIM. It is much more likely that I am a simulant person than a non-simulant person.¹

¹ Along with ‘might’, above, ‘likely’ here has an epistemic sense: likelihood on my evidence. So SIM could be true even if, unknown to me, conscious simulations are metaphysically impossible. I explain in more detail how I understand SIM and related claims in section 2.

✉ Teruji Thomas
teruji.thomas@philosophy.ox.ac.uk

¹ Global Priorities Institute, Faculty of Philosophy, University of Oxford, Oxford, UK

Bostrom (2003) presents a closely related argument (with a correction in Bostrom and Kulczycki (2011)), known as the Simulation Argument. It has inspired a great deal of philosophical and popular discussion. Strictly speaking, the Simulation Argument, as Bostrom construes it, is *not* an argument for SIM, but it does *suggest* such an argument, which has been the focus for most discussion. I'll henceforth call this suggested argument 'the classic argument'—remembering that Bostrom himself does not endorse it. Arguments similar to the classic one have been given for other hypotheses, for example the hypothesis that I am more likely to be a Boltzmann brain than an ordinary person²—but I'll focus on simulations here.

I'll give a reconstruction of the classic argument in section 2, but I can already explain why I find it unsatisfactory. The argument depends on the premiss that my current evidence entails (or at least strongly supports) the following empirical claim:

HIGH RATIO. The ratio of simulant to non-simulant people is high.

Bostrom did not endorse HIGH RATIO, and he appears to end up with roughly a 1/3 credence in it. For all the classic argument says, this is compatible with a 1/3 credence that I am a simulant person and a 2/3 credence that I am a non-simulant person. (This is why Bostrom does not endorse SIM, although of course a 1/3 credence is still surprisingly high.) More generally, for all the classic argument says, it could be *much* less likely that I'm a simulant person than a non-simulant person, as long as HIGH RATIO is unlikely.

Thus, one problem for the classic argument is that I don't seem to have much evidence for HIGH RATIO. Perhaps more interestingly, it seems hard to imagine evidence that *would* straightforwardly support HIGH RATIO, unless I can rely on the assumption that I'm a non-simulant person—an assumption the argument claims is unlikely to be true!³ To take the simplest example, suppose I discovered that the number of people in my world was much smaller than the number of simulant people they create. This would only straightforwardly support HIGH RATIO on the assumption that I live in the ground-level of reality, so that the people in my world are non-simulant people. More generally, if I am in a simulation, then I have no experience of the ground-level of reality, and so little information about how many non-simulant people there are, or how many simulations they create. I can gather evidence about the world I live in, but it is not clear why the ground-level of reality must be similar to my own.

² See Carroll (2020); Kotzen (2020); Dogramaci (2020) for recent discussion. In this paper, I'm happy *not* to count Boltzmann brains and other 'freak observers' (Crawford, 2013) as 'people'. Indeed, I'll sometimes assume that, if I'm a non-simulant person, then my experiences are generally veridical, and the world is basically as we ordinarily take it to be. On the other hand, this means that (i) it may not be certain on my evidence that I am a person at all; (ii) SIM does not say anything about whether I am more likely to be a simulant person or a non-simulant freak observer—a question I simply set aside.

³ This brief line of thought lies at the heart of objections that the classic argument is self-undermining (SIM undermines whatever reasons I might have had to believe HIGH RATIO), that it relies on 'selective skepticism' (Birch, 2013, esp. section 3), or that it results in 'epistemic instability' (Crawford, 2013); see Dogramaci (2020) for similar objections in the context of Boltzmann brains. One can thus read this paper as a response to theirs: I give an argument for SIM on roughly the same grounds but immune to these worries. Brueckner (2008) objects more simply that the probability of HIGH RATIO is inscrutable; this does not affect my argument either.

So my own argument will not be based on HIGH RATIO, and this is the main way in which it improves upon the classic argument. Instead, the analogous premiss in my argument is

HIGH EXPECTATION. *Conditional* on my being a non-simulant person, the *expected* ratio of simulant to non-simulant people *in the appropriate reference class* is high.

Note the three key ways, marked by italics, in which HIGH EXPECTATION differs from HIGH RATIO. First, it involves the condition that I am a non-simulant person. When combined with this condition, my ordinary empirical evidence *is* informative about the ground-level of reality, and may be used to support claims about the ratio of simulant to non-simulant people. Second, HIGH EXPECTATION may be true even if this ratio is unlikely (on that condition) to be high; the expected ratio may be high, even if the actual ratio is almost certainly zero. Finally, HIGH EXPECTATION concerns not all people, but only those ‘in the appropriate reference class’. I’ll say more about this restriction later, but the basic idea is to consider only people who inhabit worlds in broad strokes like our own: they live on minor variants of 21st century Earth.⁴

Although I will not argue for HIGH EXPECTATION in any detail, I believe it is *this* claim, and not HIGH RATIO, that is supported by the empirical considerations adduced in Bostrom’s paper, and especially by his claims about feasible computing power. Here’s the idea. Suppose I’m a non-simulant person. It may be quite unlikely that our descendants will run simulations of their ancestral 21st century. But (Bostrom argues) they might in principle run *enormously* many, at negligible cost to themselves, and even on a whim.⁵ Those simulated 21st century people would be in the specified reference class. So there’s at least a small probability that the ratio of simulant to non-simulant people in the reference class is enormous. We can use ordinary empirical considerations to estimate this probability. As long as the probability is not *too* small, HIGH EXPECTATION is true. In contrast, this line of reasoning does not seem to support HIGH RATIO.

The main advance in this paper is to replace HIGH RATIO by HIGH EXPECTATION, thus giving us more reason to take seriously the possibility that we are simulant people. I develop the argument in sections 2 to 3, and then discuss some remaining issues in section 4. I will especially consider, and tentatively respond to, a worry raised by Weatherston (2003) about what evidence I can have if I happen to be in the ground-level of reality. And I will refute the common line of thought that finding lots of simulations being run must increase the odds that I myself am in one.

Let me conclude this introduction by explaining why I think assessing SIM is an important project, even from a practical point of view. Jenkins (2006, p. 23) puts the

⁴ Of course, we could have included a similar restriction in HIGH RATIO, or just incorporated it into the definition of ‘people’. But it has not been relevant until now.

⁵ According to Bostrom (2003, pp. 247–8), ‘A single [planetary-mass] computer could simulate the entire mental history of humankind...by using less than one millionth of its processing power for one second. A posthuman civilization may eventually build an astronomical number of such computers.’

point in a usefully direct and provocative way. Given the classic argument, he writes, it is

highly probable that we are a form of artificial intelligence inhabiting one of these simulations. To avoid stacking (i.e. simulations within simulations), the termination of these simulations is likely to be the point in history when the technology to create them first became widely available.... Long range planning beyond this date would therefore be futile.

While this line of thought is highly speculative, it does seem right that living in a simulation may carry with it some distinctive risks (e.g., the risk of simulation termination) as well some distinctive opportunities (e.g., to mitigate termination risk by restricting technological progress).⁶ These are, in other words, speculations that may be worth taking seriously if there really is a good argument for SIM. My goal here is to give the best and clearest argument I can.

2 Preliminaries

In this section, I'll explain the framework I'll be using in this paper, and then I'll indicate the general form of the argument to follow—including a brief reconstruction of the classic argument for SIM.

2.1 The Framework

First, the basic framework. Questions of likelihood could be interpreted in different ways. I will use a Bayesian framework, in which the natural question is what is supported, probabilistically, by my total evidence. So, SIM is the claim that my current evidence strongly supports the hypothesis that I'm a simulant person over the hypothesis that I'm a non-simulant person.

As usual in this Bayesian framework, my total evidence and the hypotheses I entertain are all *propositions*. Following Lewis (1979), I think of each proposition as attributing some property to myself. For example, the proposition that grass is green is equivalent, for my purposes, to the proposition that I'm in a possible world where grass is green. For any property F , let iF be the proposition that I am F . So if F is the property of being a simulant person, then iF is the proposition that I am a simulant person. In the same way, let iE denote my total evidence, whatever proposition it may be.

Instead of just tracking my current credences, I assume that the facts (or, if you prefer, my subjective judgments) about evidential support are encoded by a probability measure (or Popper function) Pr , my *ur prior*. The probability that I am F , given my total evidence that I am E , is thus the conditional probability $\text{Pr}(iF \mid iE)$.

⁶ For more recent and careful work exploring these and similar worries, see Tomasik (2016) and especially Greene (2020). See Chalmers (2022) for a wider-ranging exploration of life in virtual worlds.

That said, it is sometimes more convenient to talk in terms of odds rather than probabilities. (The specific reason it is convenient will appear at the end of §3.3.) The odds that I am F rather than G , given that I am E , are defined to be

$$\text{Odds } ({}_tF/{}_tG \mid {}_tE) \stackrel{\text{def}}{=} \frac{\Pr({}_tF \mid {}_tE)}{\Pr({}_tG \mid {}_tE)}.$$

So SIM says that the odds are high that I'm a simulant person (${}_tF$) rather than a non-simulant person (${}_tG$), given my current evidence (${}_tE$).⁷

2.2 The Classic Argument

With this framework in place, the general style of my argument—and of the classic argument—can be anticipated from the following highly plausible, if informal, piece of reasoning.

I'm a smoker, and one out of fifteen smokers develops lung cancer. So, if I had no other information, I would say that the probability that I'll develop lung cancer is 1/15. I do have other information, but none of it is particularly relevant to whether I'll develop lung cancer. So, taking all my evidence into account, the probability is still about 1/15.

This reasoning depends on (1) the identification of a 'reference class'—the class of smokers; (2) an empirical premiss about frequencies within that reference class—the frequency of lung cancer; (3) a premiss about what could be concluded if one had no other information; (4) the premiss that, although one does have other information, it is mostly irrelevant to the question at hand.

The classic argument follows this logic quite closely.⁸ (1) For the reference class, we can take the class of all people. (2) The empirical premiss is HIGH RATIO, to the effect that the vast majority of this reference class are simulants. (3) If I had no other information, this would make it extremely likely that, if I am a person at all, then I am a simulant person. (4) Finally (Bostrom suggests) though I do have other evidence, none of it tells very strongly *against* the thesis that I am a simulant person rather than a non-simulant person. Therefore, even taking all my evidence into account, it is still far more likely that I am a simulant person than a non-simulant person.

In the introduction, I raised some doubts related to premiss (2). I also think that premiss (4), while somewhat plausible, is difficult to assess. I will discuss some worries about it in section 4, focusing (however) on the analogous step in my own argument for SIM. For now, let me make some comments about steps (1) and (3) in this general form of argument; I will otherwise assume they are unproblematic.

⁷ As mentioned in footnote 2, I don't assume that it's certain that I'm a person; thus *I'm not a simulant person*, $\neg {}_tF$, may not be equivalent to *I'm a non-simulant person*, ${}_tG$. If one doesn't care about this issue, one can replace G by not- F in all the relevant places.

⁸ The form of the argument reconstructed here is more overt in the exchange between Weatherson (2003) and Bostrom (2005) than in the original paper (Bostrom, 2003).

When it comes to step (1), the specification of the reference class is entirely a matter of stipulation: it can be chosen once and for all in whatever way is convenient for the argument at hand. The key points are that I have relevant information about frequencies within the reference class, as required by step 2; and that I don't have much relevant evidence beyond this, as required by step 4. (There can be some tension between these requirements, as I discuss in section 4.) My use of reference classes is thus different from the way they are used in some of the neighbouring literature on anthropic reasoning, where identifying the *correct* reference class is a vexed topic, with different reference classes leading to contradictory results (Bostrom, 2002).

As for step (3), I will state a precise version of the principle all these arguments appeal to in section 3.3, where I call it CALIBRATION. But it is already worth clarifying that it is not a completely general (and therefore implausible!) principle about aligning credences with observed frequencies. It is, instead, a claim specifically about *self-location*: insofar as my evidence is compatible with different hypotheses about who I am within a given world, these hypotheses are equally likely. To return to our initial example, if literally all I know about myself is that I am a smoker, then (the principle claims) it is equally likely that I am any given smoker within each world. If I also know that one in fifteen smokers develops cancer, then (it follows) the probability that I'm one of *those* smokers is 1/15. This type of principle is relatively uncontroversial in the literature on self-locating credences.⁹ At any rate, it is common ground between my argument and Bostrom's, so I intend to rest upon its intuitive plausibility here.

3 The Argument

The structure of the argument will parallel those described in section 2.

3.1 The Reference Class

As I mentioned in the introduction, I take the reference class to consist of people living on minor variants of 21st century Earth. At least, this is a provisional definition that is enough to get the argument going. We can amend it in whatever way makes the premisses most compelling—a point on which I'll elaborate as we go, especially in section 4.

⁹ Bostrom (2003, §5) calls it a 'bland indifference principle'. In the language of centered worlds, the claim is that different centerings of the same world are equally likely, insofar as they are compatible with my evidence. Elga (2004) gives positive arguments for this view, which he calls 'INDIFFERENCE'. However, see Weatherson (2005) for a dissenting view, according to which the probabilities of self-locating hypotheses are often imprecise. And see footnote 11 for some further comments.

3.2 The Empirical Premiss

The first premiss of the argument is HIGH EXPECTATION, which is basically an empirical claim about the ratio of simulant to non-simulant people in the reference class. (More exactly, it is a claim about my evidence about this ratio.)

Let me give a more formal statement of the premiss. Let F be the property of being a simulant person, G the property of being a non-simulant person, and R the property of being in the reference class. (For short I will sometimes refer to R itself as ‘the reference class’.) Let $\text{Rat}_R^{F/G}$ denote the ratio of F s to G s among all R s (i.e., here, the ratio of simulant to non-simulant people in the reference class). Then the expected ratio of simulant to non-simulant people in the reference class, conditional on my being a non-simulant person, is defined by the sum

$$\mathbb{E}(\text{Rat}_R^{F/G} \mid \iota E \ \& \ \iota G) \stackrel{\text{def}}{=} \sum_r r \times \Pr(\text{Rat}_R^{F/G} = r \mid \iota E \ \& \ \iota G).$$

Here r ranges over the countably many candidate values for $\text{Rat}_R^{F/G}$, i.e., all non-negative rational numbers.¹⁰ With this definition, the premiss is

HIGH EXPECTATION. $\mathbb{E}(\text{Rat}_R^{F/G} \mid \iota E \ \& \ \iota G)$ is high.

I already sketched in the introduction why HIGH EXPECTATION might be true with respect to the specified reference class. Of course, ‘high’ is vague and context-dependent, but this doesn’t matter, since the argument will conclude (in equation 1 below) that the odds I’m a simulant person are also high in whatever sense is on the table.

3.3 Frequency-Based Reasoning

Let’s revisit the initial example from section 2. Suppose that (strangely enough) my only evidence is that I’m a smoker and that 1 in 15 smokers develop lung cancer. How likely is it, on this evidence, that I develop lung cancer? As I suggested, the probability is 1/15, calibrated to the frequency of cancer-developers among all smokers. Or, in other words, the odds are 1/14, calibrated to the ratio of cancer-developers to cancer-non-developers among all smokers.

¹⁰ A non-trivial assumption I’m making is that the number of R s is finite. Then $\text{Rat}_R^{F/G}$ is also finite, on the condition that I myself am both R and G . In universes with infinite populations, frequency-based reasoning faces general and serious problems (Arntzenius & Dorr, 2017)

Here's the obvious generalization, formulated for any properties F , G , and R , and any non-negative rational number r .¹¹

CALIBRATION. Odds ($tF/tG \mid tR$ & $\text{Rat}_R^{F/G} = r$) = r .

In words: given only that I'm an R and that the ratio of F s to G s among R s is r , the odds are r that I'm an F rather than a G . The second premiss of my argument is that CALIBRATION holds for all r , when F is the property of being a simulant person, G is the property of being a non-simulant person, and R is the reference class.

In the appendix (and illustrated by the example below) I prove that CALIBRATION has the following non-obvious formal consequence, which is what's directly relevant to my argument:

CALIBRATION*. Odds ($tF/tG \mid tR$) \geq $\mathbb{E}(\text{Rat}_R^{F/G} \mid tR \ \& \ tG)$.

So, if my total evidence consisted in the facts that characterise the reference class (tR), then the odds that I am a simulant person rather than a non-simulant person would be at least as great as the expected ratio of simulant to non-simulant people, conditional on my being a non-simulant person.

By the way, CALIBRATION* explains why I generally focus on ratios and odds. The principle is more cumbersome when expressed in terms of frequencies and probabilities.

3.4 Admissibility

The final premiss claims that I do not have much relevant evidence beyond the fact that I am in the reference class. We can include in the proposition tR facts about the laws of physics, the limits of computational power, human psychology, and the apparent trajectory of civilization (all, that is, for the world I inhabit, not necessarily for the ground-level of reality). And I simply don't have much to go on, beyond these general facts, when it comes to assessing (i) whether I am a simulant person or not, or (ii) the ratio of simulant to non-simulant people among the R s. The more specific details of my life, like my name and address, certainly don't seem relevant.¹²

¹¹ As I remarked at the end of section 2, CALIBRATION can be understood as a principle of indifference between different self-locating hypotheses. Let me mention two subtleties that I consider more carefully in Thomas (2021). First, it's possible that I am F at some times and not- F at others, so frequencies have to be understood in a time-weighted way. This is related to the point that I can be uncertain not only who I am, but when I am. Second, problems arise unless F , G , and R are presented in what Chalmers (2004) calls 'neutral' terms. To give a trivial example, suppose that my total evidence is that I am a person, and that one in 8 billion people is me. I can still be certain that I am me! In that paper I defend a more elaborate principle, which arguably reduces to CALIBRATION in the case at hand.

¹² We could always change tR to include any further facts that *do* seem relevant, but this may weaken the argument for HIGH EXPECTATION—a tradeoff that I will discuss more in section 4. One idea that appears in the literature is that there might be disproportionately many simulations of *important* or otherwise *interesting* lives. Doesn't the fact that *my* life is uninteresting provide further evidence that I am not in a simulation? The main response, following Bostrom (2005, p. 93), is that any such considerations are 'at best weak and speculative'. But it is also worth noting that my uninteresting life may be part of a whole-world simulation that is motivated by the interesting lives of other people. The 'uninteresting life' consideration speaks most plausibly against the hypothesis that I am in a 'selective' (Bostrom, 2003, p. 254) or 'sol-

At first, let us suppose that my additional evidence is *completely* irrelevant in the two ways just mentioned. Formally, then, the reference class R satisfies the following two conditions (in which case I'll say that R is *admissible*).

ADMISSIBILITY.

- (i) $\text{Odds}(tF/tG \mid tE) = \text{Odds}(tF/tG \mid tR)$;
- (ii) $\mathbb{E}(\text{Rat}_R^{F/G} \mid tE \ \& \ tG) = \mathbb{E}(\text{Rat}_R^{F/G} \mid tR \ \& \ tG)$.

The point of ADMISSIBILITY is that it allows us transform each side of the inequality in CALIBRATION* to get:

$$\text{Odds}(tF/tG \mid tE) \geq \mathbb{E}(\text{Rat}_R^{F/G} \mid tE \ \& \ tG). \quad (1)$$

This is like CALIBRATION*, except that it takes my total evidence (tE) into account. At this point we can appeal to HIGH EXPECTATION. It tells us that the expected ratio is high; so, therefore, are the odds.

Now, more plausibly, conditions (i) and (ii) will only be approximately true, so that R is only approximately admissible. But this does not make much difference. Indeed, as long as (i) $\text{Odds}(tF/tG \mid tE)$ is not much smaller than $\text{Odds}(tF/tG \mid tR)$, and (ii) $\mathbb{E}(\text{Rat}_R^{F/G} \mid tE \ \& \ tG)$ is not much larger than $\mathbb{E}(\text{Rat}_R^{F/G} \mid tR \ \& \ tG)$, it follows from CALIBRATION* that $\text{Odds}(tF/tG \mid tE)$ is not much smaller than $\mathbb{E}(\text{Rat}_R^{F/G} \mid tE \ \& \ tG)$. If the latter is high, so must be the former.

This concludes the argument for SIM, except for the proof of CALIBRATION*, which I give in the appendix. In the rest of this section, I'll give a 'proof by example'—that is, a toy example that illustrates the logic of the proof, which many readers may find more illuminating than the derivation in full generality.

3.5 Example

Suppose that the entirety of my evidence is that I am R —it is this situation that CALIBRATION* concerns. To get a simple example, let me stipulate the following details. There are two possibilities, A and B , compatible with the assumption that I am both R and a non-simulant person. Either

- (A) There are no simulant people who are R s, so $\text{Rat}_R^{F/G} = 0$; or else

Footnote 12 (continued)

ipish' (Tomasik, 2016) simulation that contains only a few people (perhaps just me!) with rich mental lives, plus a bunch of mindless zombies. At the beginning of the paper I framed the topic in terms of 'simulations of whole worlds'; if we take that framing seriously, then the possibility that I am in a selective simulation—however unlikely it may be—is an *additional* way that I might fail to live in the ground-level of reality. (My discussion of 'vast and mostly lifeless worlds' in section 4 is also relevant here.)

- (B) There are 100 simulant people for every non-simulant person among the R s, so $\text{Rat}_R^{F/G} = 100$.

Suppose more specifically that

- (*) Conditional on my being a non-simulant person, A has probability 0.9 and B has probability 0.1.

It follows that the expected ratio of simulant to non-simulant people is

$$\begin{aligned}\mathbb{E}(\text{Rat}_R^{F/G} \mid \iota R \ \& \ \iota G) &= 0 \times \text{Pr}(A \mid \iota R \ \& \ \iota G) + 100 \times \text{Pr}(B \mid \iota R \ \& \ \iota G) \\ &= 0 \times 0.9 + 100 \times 0.1 \\ &= 10.\end{aligned}$$

Now we have to calculate the odds that I'm a simulant person. Let's start by comparing the probabilities of the following three propositions, conditional on my total evidence, ιR .

($A \ \& \ \iota G$) A is true and I'm a non-simulant person.

($B \ \& \ \iota G$) B is true and I'm a non-simulant person.

($B \ \& \ \iota F$) B is true and I'm a simulant person.

First, $A \ \& \ \iota G$ is nine times more likely than $B \ \& \ \iota G$. How come? We can use the mathematical identity $\text{Odds}(A \ \& \ \iota G / B \ \& \ \iota G \mid \iota R) = \text{Odds}(A / B \mid \iota R \ \& \ \iota G)$. The latter equals 9, according to (*).

Second, $B \ \& \ \iota F$ is 100 times more likely than $B \ \& \ \iota G$. To get this, we can use the mathematical identity $\text{Odds}(B \ \& \ \iota F / B \ \& \ \iota G \mid \iota R) = \text{Odds}(\iota F / \iota G \mid \iota R \ \& \ B)$. The latter equals 100, by CALIBRATION.

It follows from the previous two claims that $B \ \& \ \iota F$ is 10 times more likely than $A \ \& \ \iota G$ and $B \ \& \ \iota G$ put together.¹³ But $A \ \& \ \iota G$ and $B \ \& \ \iota G$ exhaust the possibilities for my being a non-simulant person, while $B \ \& \ \iota F$ does not necessarily exhaust the possibilities for my being a simulant person. So it's *at least* 10 times more likely than I'm a simulant person than a non-simulant person. We have established

$$\text{Odds}(\iota F / \iota G \mid \iota R) \geq 10 = \mathbb{E}(\text{Rat}_R^{F/G} \mid \iota R \ \& \ \iota G).$$

Thus we have derived the key claim, CALIBRATION*, in this special case.

¹³ In detail: since A and B are mutually incompatible, we have $\text{Pr}((A \ \& \ \iota G) \vee (B \ \& \ \iota G) \mid \iota R) = \text{Pr}(A \ \& \ \iota G \mid \iota R) + \text{Pr}(B \ \& \ \iota G \mid \iota R)$, which equals $10 \text{Pr}(B \ \& \ \iota G \mid \iota R)$ according to the first claim. The second claim says $\text{Pr}(B \ \& \ \iota F \mid \iota R) = 100 \text{Pr}(B \ \& \ \iota G \mid \iota R)$. So, in combination, $\text{Pr}(B \ \& \ \iota F \mid \iota R) = 10 \text{Pr}((A \ \& \ \iota G) \vee (B \ \& \ \iota G) \mid \iota R)$.

4 Discussion

The main subtlety in my argument for *SIM* concerns the choice of the reference class *R*. I have so far used a rough characterization of *R* in terms of ‘minor variants of 21st century Earth’, which is enough to show how such an argument can get off the ground (whereas the classic argument seems hopeless). However, since the conclusion, *SIM*, does not depend on *R*, we can choose *R* in any way that makes the premisses compelling. Of these premisses, I will not discuss *CALIBRATION* further, since it reflects a very general form of frequency-based reasoning that should work for most reasonable ways of specifying *R*. But there is an interesting tension between the remaining premisses, *HIGH EXPECTATION* and *ADMISSIBILITY*, which I will consider in this section. To be clear: essentially the same tension exists in the classic argument for *SIM*, and I don’t think my own argument is much, if at all, worse off in this respect. But the following points are still important for understanding the limitations of the argument I’ve proposed.

The basic tension is as follows. We can guarantee that *ADMISSIBILITY* holds by packing a lot of information into *R*, making the proposition ιR close to my total evidence; but this makes *HIGH EXPECTATION* less plausible. For example, *ADMISSIBILITY* holds trivially if we take $R = E$, for then I have no further evidence at all beyond ιR . But, supposing that I’m a non-simulant person, it’s unlikely that there are many, or perhaps any, simulant people with the very specific property *E*. So the expected ratio will be relatively low. On the other hand, if we take *R* to be much weaker than my total evidence—applying, perhaps, to most people who have ever lived—then *HIGH EXPECTATION* may be arguable, but *ADMISSIBILITY* becomes less plausible, or simply hard to assess. I will study two examples of this phenomenon.

Before proceeding, let me make explicit a point that could easily get lost in the high-level conceptual discussion: the actual numbers matter! If the expected ratio in *HIGH EXPECTATION* is extremely high, then the odds I’m a simulant will also be high, unless *ADMISSIBILITY* dramatically fails.

4.1 The Limits of Appearances

Let’s suppose that I am in fact a non-simulant person, living in the ground-level of reality. And for concreteness let’s consider the view that my evidence consists of what I know (Williamson, 1997). One thing I know—or so we would ordinarily say—is that the world around me is more than 10 billion years old, and that throughout many of those billions of years, vast expanses of the universe were lifeless. In contrast, we might think that most simulated worlds are *not* vast and mostly lifeless (*vML*) in this way: it would be easier to simulate a relatively small, inhabited portion of the universe and make it *appear* that the world was *vML*.

This leads to a dilemma: should we include in the reference class property *R* the fact that the world is *vML*? The first horn: if we do include it, then the expected ratio of simulant to non-simulant *R*s will be relatively low, since most simulant worlds are *not* *vML*. The second horn: if we don’t include it, then the first clause of *ADMISSIBILITY*

will fail, since I have additional evidence (the fact that my world is v_{ML}) that confirms I'm a non-simulant.¹⁴

I see this dilemma as the philosophically most interesting objection to my argument, so let me take some time to sketch, non-committally, how it might be defused. One could, of course, try to argue on empirical grounds that HIGH EXPECTATION is true even if we include the claim that the world is v_{ML} , and similar claims, in R . But otherwise I see three promising, if closely related, conceptual strategies.

The first is to adopt an internalist conception of evidence, according to which my evidence is essentially just a matter of how things appear to me. So, even if I'm a non-simulant person, my evidence does not entail that my world is v_{ML} ; perhaps it only entails that my world *appears* v_{ML} . There is no problem with including this fact about appearances in the reference class, insofar as it would not greatly reduce the expected ratio.¹⁵

Alternatively, we could retain an externalist conception of evidence (as I myself would prefer), and even retain the view that my evidence is my knowledge, but argue that there are special features of the case that prevent me from knowing that my world is v_{ML} . Bostrom (2005, p. 94) suggests a view along these lines, in the context of the classic argument for SIM. In a possible world where there is a high ratio of simulant to non-simulant people, 'illusions are ubiquitous': 'almost all people...have perceptions which, if interpreted naïvely, are misleading about [certain] facts'. And if we know that's how things are, then we cannot trust the appearances: knowledge of HIGH RATIO (or maybe the mere truth of HIGH RATIO) prevents me from knowing that my world is vast and mostly lifeless.

Be that as it may, it is initially unclear whether Bostrom's move helps defend my version of the argument, the whole point of which is not to rely on HIGH RATIO. I'm not considering a situation in which I *know* that the ratio of simulant to non-simulant people in my reference class is high, or even necessarily one in which *it's true* that this ratio is high. It's just that the *expected* ratio is high, conditional on my being a non-simulant person. To escape the dilemma in Bostrom's way, I would need to argue that in *this* epistemic situation I cannot know that my environment is v_{ML} .

¹⁴ The claim that my world is v_{ML} is, of course, only an example. Weatherston (2003, pp. 429–430) raises a similar objection in the context of the classic argument for SIM, though he focuses on evidence about more proximate matters like the claim that I have eyeballs. Bostrom (2005) and especially Chalmers (2005) provide some ammunition to respond to Weatherston: they argue that, if I am in a sufficiently rich simulation, then I *do* have eyeballs and the like. Indeed, whether or not I'm in such a simulation, my word 'eyeballs' picks out the sort of visual apparatus that people in my world have, so the fact that I have eyeballs doesn't provide me with evidence either way. In general, being in a simulation need not be the sort of radically skeptical scenario in which my most commonplace beliefs turn out to be false. I focus on the claim that my world is v_{ML} because, as Chalmers acknowledges in the various scenarios considered in his section 8, the world beyond my immediate macroscopic environs is something about which I, as a simulant, could more easily be mistaken.

¹⁵ Weatherston (2003, p. 430) suggests that a version of the objection applies to some internalists as well. To place some limits on the discussion, I have to set aside the interesting question of what the best internalist views, and other externalist views, would say here; I don't want to suggest that they go scot free. In particular, it is important to remember that the basic Bayesian framework requires one's evidence to be a proposition, and on this score it is not always clear what internalists have in mind.

While I lack a compelling argument for this position, I think it is defensible. Let me explain. Suppose we include in the specification of the reference class only claims about how the world appears, such as the claim that the world appears to be v_{ML} . The dilemma primarily concerns the first clause of *ADMISSIBILITY*, not the second, which seems relatively safe.¹⁶ Using only that second clause, *CALIBRATION** still leads to the inequality

$$\text{Odds}(I F / I G \mid I R) \geq \mathbb{E}(\text{Rat}_R^{F/G} \mid I E \ \& \ I G).$$

Continuing to assume that the expected ratio on the right-hand side is high, this inequality shows that the general facts about how things appear (those involved in the fact that I am R) strongly support the view that I am a simulant rather than a non-simulant person. Let us also grant to the objector that no simulated worlds are really v_{ML} . Then, going by these general facts about how things appear—including the fact that my world appears to be v_{ML} —it's highly unlikely that my world is v_{ML} . This does seem like the type of situation that could prevent my knowing that my world is v_{ML} .¹⁷

The third strategy for avoiding the dilemma is a variation on the second. We can make a partial retreat. There are different questions to be asked: questions about what I'm in a position to know and what that knowledge would support; but also questions about what credences a reasonable but fallible subject would have or ought to have.¹⁸ We might concede to the objector that, *if* in fact I'm a non-simulant person, then I'm in a position to know that my world is v_{ML} , and my evidence would then support high confidence that I'm a non-simulant person. But, in another sense, I'm still in an unusually bad position to affirm the antecedent: as I've just argued, the facts about how things broadly appear support its negation. This may make it reasonable to adopt credences that are out of line with what's supported by my evidence, or—perhaps better—to abandon my belief, and hence my knowledge, that the world is, in the relevant ways, exactly as it seems.

4.2 The Limits of Future Evidence

While I sketched an argument for *HIGH EXPECTATION* in §1, the success of that argument depends on empirical factors that I have not examined closely. On the other hand, even if that argument fails, it seems plausible that we could gain evidence in the future that would make *HIGH EXPECTATION* true. However, contrary to some claims in the literature, the overall effect of such evidence can be difficult to gauge, and would not necessarily tend to confirm that we are simulant people.

¹⁶ That second clause involves the condition that I am a non-simulant person, and we may assume that the appearances are generally veridical on this condition (cf. footnote 2). So, as far as the second clause goes, there is not much space between the claim that my world appears to be v_{ML} and the claim that it actually is.

¹⁷ To be more explicit: what's highly unlikely is the conditional that, if I'm a person at all, then my world is v_{ML} . It's knowledge of this conditional that is required to support the claim that I'm a non-simulant person as opposed to a simulant person.

¹⁸ Without meaning to attribute to them any position on the current topic, I have in mind the kinds of distinctions drawn by Lasonen Aarnio (2010) and Schoenfeld (2012).

For example, what would happen if I were to find a secret lab running quadrillions of whole-world simulations, or if I were to find myself about to run such simulations? Bostrom (2006, p. 9) and Greene (2020) both claim that this should make me confident that I am a simulant person.¹⁹ In terms of my own argument, I am happy to concede that such evidence could make HIGH EXPECTATION true. The problem is that the reference class may no longer be admissible.

Before explaining that problem in detail, I think it is worth taking a step back. Instead of applying a general argument for SIM, we can just ask what kind of an update the new evidence would provide, using updating norms like conditionalization. The answer in part depends on what we take the new evidence to be. On the externalist account, perhaps I get to know that there *really are* many simulations being run within my environment. But as Lewis (2013) explains, this evidence *could* quite powerfully confirm that I myself am non-simulant, since it's plausible on combinatorial grounds that the vast majority of simulations do not themselves contain many simulations.²⁰ Even if we focus on how things appear, the situation is unclear. Perhaps the new evidence is only that there *appear* to be many simulations within my environment. (After all, the lab computers bleep and bloop, but I don't have direct access to what's going on inside.) In order to strongly confirm that I'm a simulant person, this appearance must be much more likely on the supposition that I'm a simulant than on the supposition that I'm a non-simulant, against the background of my current evidence.²¹ I don't see why it would be.

Returning to the admissibility constraint, there are two ways to look at the situation. If we keep the reference class fixed, then my new evidence may well make the expected ratio high. But insofar as this new evidence supports the claim that I am a simulant person, the first clause of ADMISSIBILITY must fail. On the other hand, we could analyse things in terms of a new reference class. We must choose the reference

¹⁹ As Greene discusses at length, this claim could be practically relevant to whether we *should* run simulations, at least given evidential decision theory. Bostrom and Greene are both thinking about the classic argument, and so for them the idea is that finding the secret lab would make me confident in HIGH RATIO. That's hardly obvious: again, if I'm in a simulation, finding a secret lab would not tell me much about the number of non-simulant people. But even granting this claim, step (4) in the classic argument might well fail with respect to my new evidence, for reasons similar to those I discuss below.

²⁰ To illustrate, suppose that people in the ground-level world created 100 simulated worlds, each of which contained a further 100 simulated worlds (and the process stopped there). Then 99% of worlds would contain no simulations. Curiously, Lewis seems to construe the key point of the Simulation Argument as being a probabilistic inference from *there are no simulations in my world* to *I'm in a simulation*.

²¹ Here I'm appealing to the fact that when my evidence strengthens from tE to tE' , the odds that I am F rather than G get multiplied by the *Bayes factor*

$$\text{Bayes}_{tF/tG}(tE' | tE) = \frac{\Pr(tE' | tF \& tE)}{\Pr(tE' | tG \& tE)}.$$

In words: the odds increase insofar as it is more likely that I am E' on the supposition that I am F than on the supposition that I am G , against the background evidence that I am E . It may be worth reiterating the point made by Crawford (2013, p. 262) that finding lots of (real or apparent) simulations would hardly make it likely that I'm living in one of *those* simulations. The situation may become complicated if the discovered simulations involve phenomenal duplicates of me, leaving me in the position of 'Dr. Evil' (Elga, 2004). Alternatively, perhaps there could be metaphysical cycles of simulations within simulations—thanks to Kenny Easwaran for raising this point.

class property R so that the fact that I am R includes most of my relevant evidence. Obviously, one relevant part of my evidence is that there are (or there appear to be) many simulations run within my world. So we have to include something to that effect within the definition of R —we have to consider a more limited reference class. At that point, the problem is that the expected ratio with respect to this new reference class need not be higher than before.

5 Conclusion

The situation is this. Bostrom's Simulation Argument suggests a further argument, based on HIGH RATIO, that we ourselves are simulants. But it's hard to see why we should be confident of HIGH RATIO, and it's hard to see how we could get good evidence about the ratio of simulant people to non-simulant people if we ourselves are in a simulation. This paper indicates how to get around both these points. First, it suffices to have evidence about the ratio *on the condition* that we are non-simulant people. Second, we need not be confident, on that condition, that the ratio is high; it suffices that the *expected* ratio is high. Because of this, the resulting argument for SIM is much more troubling (if that is the right word) than the one based on HIGH RATIO.

That argument is the main point of this paper. However, I have also indicated some ways to resist it. A key step in the argument is to assume that there is some reference class property R that is both roughly admissible and that leads to a high expected ratio. One could resist this on strictly empirical grounds (for example, though not only, by rejecting Bostrom's claims about the limits of computational power or the uses to which our descendants seem likely to put it). Or perhaps one could resist this step on more conceptual grounds by adopting an appropriate theory of evidence. Tentatively, though, I have suggested that giving high odds that we are simulant people may still be reasonable *even if* we are non-simulant people and our evidence supports the opposite conclusion. Finally, I have argued that finding a large number of simulations being run, or otherwise acquiring evidence that raises the expected ratio, would not necessarily raise the odds that we are simulant beings. All of these topics deserve further consideration.

Appendix: Proof of Calibration*

The main technical step in the argument for SIM is the derivation of

$$\text{CALIBRATION}^*. \text{ Odds}(tF/tG \mid tR) \geq \mathbb{E}(\text{Rat}_R^{F/G} \mid tR \ \& \ tG).$$

This only requires CALIBRATION. I illustrated the derivation in section 3 with a toy example, but here I consider the general case.

Start from the observation that

$$\Pr(tF \mid tR) \geq \sum_r \Pr(tF \ \& \ \text{Rat}_R^{F/G} = r \mid tR) \quad (2)$$

where the sum is over all non-negative rational r . I will consider each term in this sum, one at a time. Let $a(r)$ denote the r th term:

$$a(r) \stackrel{\text{def}}{=} \Pr(\iota F \ \& \ \text{Rat}_R^{F/G} = r \mid \iota R).$$

Recall the definition of conditional probability: $\Pr(X \mid Y) = \Pr(X \ \& \ Y) / \Pr(Y)$. I will use this to repackage conditional probabilities without detailed explanation. (If one prefers to treat conditional probabilities as primitive, one could instead use the axioms for Popper functions.) As a first application, for each candidate value of r we have

$$a(r) = \Pr(\iota F \mid \iota R \ \& \ \text{Rat}_R^{F/G} = r) \times \Pr(\text{Rat}_R^{F/G} = r \mid \iota R).$$

Now we can express the first factor on the right-hand side in terms of odds and then apply CALIBRATION:

$$\begin{aligned} \Pr(\iota F \mid \iota R \ \& \ \text{Rat}_R^{F/G} = r) &= \text{Odds}(\iota F / \iota G \mid \iota R \ \& \ \text{Rat}_R^{F/G} = r) \\ &\quad \times \Pr(\iota G \mid \iota R \ \& \ \text{Rat}_R^{F/G} = r) \\ &= r \times \Pr(\iota G \mid \iota R \ \& \ \text{Rat}_R^{F/G} = r). \end{aligned}$$

Substituting this into the previous equation, we have

$$a(r) = r \times \Pr(\iota G \mid \iota R \ \& \ \text{Rat}_R^{F/G} = r) \times \Pr(\text{Rat}_R^{F/G} = r \mid \iota R).$$

Reuse the definition of conditional probability to repackage this:

$$a(r) = r \times \Pr(\text{Rat}_R^{F/G} = r \mid \iota R \ \& \ \iota G) \times \Pr(\iota G \mid \iota R).$$

If, as in the inequality (2), we sum $a(r)$ over all values of r , we get

$$\sum_r a(r) = \mathbb{E}(\text{Rat}_R^{F/G} \mid \iota R \ \& \ \iota G) \times \Pr(\iota G \mid \iota R).$$

Substituting this into (2), we get

$$\Pr(\iota F \mid \iota R) \geq \mathbb{E}(\text{Rat}_R^{F/G} \mid \iota R \ \& \ \iota G) \times \Pr(\iota G \mid \iota R).$$

Dividing through by $\Pr(\iota G \mid \iota R)$ yields CALIBRATION*.

Acknowledgements In addition to my colleagues at GPI, I'm grateful to Joe Carlsmith, Kenny Easwaran, Ben Garfinkel, Riley Harris, Maria Lasonen-Aarnio, Harvey Lederman, and Jeff Russell for helpful discussions, and to John Mori, Natasha Oughton, and Elliott Thornley for assistance.

Funding The authors did not receive support from any organization for the submitted work.

Data availability This is a purely theoretical study and there is no relevant data, code, or other materials.

Declarations

Conflict of interest There are no Conflict of interest associated with this paper.

Ethical approval This is a purely theoretical study that did not involve any human or animal participants, and no ethical approval is required.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Arntzenius, F., & Dorr, C. (2017). Self-locating priors and cosmological measures. In K. Chamcham, J. Barrow, S. Saunders, & J. Silk (Eds.), *The Philosophy of Cosmology* (pp. 396–428). Cambridge: Cambridge University Press.
- Birch, J. (2013). On the ‘simulation argument’ and selective scepticism. *Erkenntnis*, 78(1), 95–107.
- Bostrom, N. (2002). *Anthropic Bias: Observation Selection Effects in Science and Philosophy*. Routledge.
- Bostrom, N. (2003). Are we living in a computer simulation? *The Philosophical Quarterly*, 53(211), 243–255.
- Bostrom, N. (2005). The simulation argument: Reply to Weatherson. *The Philosophical Quarterly*, 55(218), 90–97.
- Bostrom, N. (2006). Do we live in a computer simulation? *New Scientist*, 192, 2579.
- Bostrom, N., & Kulczycki, M. (2011). A patch for the simulation argument. *Analysis*, 71(1), 54–61.
- Brueckner, A. (2008). The simulation argument again. *Analysis*, 68(3), 224–226.
- Carroll, S. M. (2020). Why Boltzmann brains are bad. In S. Dasgupta, B. Weslake, & R. Dotan (Eds.), *Current Controversies in Philosophy of Science* (pp. 7–20). Routledge.
- Chalmers, D. J. (2004). Epistemic two-dimensional semantics. *Philosophical Studies*, 118(1–2), 153–226.
- Chalmers, D. J. (2005). The matrix as metaphysics. In C. Grau (Ed.), *Philosophers Explore the Matrix* (pp. 132–176). Oxford University Press.
- Chalmers, D. J. (2022). *Reality+: Virtual Worlds and the Problems of Philosophy*. Norton.
- Crawford, L. (2013). Freak observers and the simulation argument. *Ratio*, 26(3), 250–264.
- Dogramaci, S. (2020). Does my total evidence support that I’m a Boltzmann brain? *Philosophical Studies*, 177(12), 3717–3723.
- Elga, A. (2004). Defeating Dr. Evil with self-locating belief. *Philosophy and Phenomenological Research*, 69(2), 383–396.
- Greene, P. (2020). The termination risks of simulation science. *Erkenntnis*, 85(2), 489–509.
- Jenkins, P. (2006). Historical simulations-motivational, ethical and legal issues. *Journal of Futures Studies*, 11(1), 23–42.
- Kotzen, M. (2020). What follows from the possibility of Boltzmann brains? In S. Dasgupta, B. Weslake, & R. Dotan (Eds.), *Current Controversies in Philosophy of Science* (pp. 21–34). Routledge.
- Lasonen Aarnio, M. (2010). Unreasonable knowledge. *Philosophical Perspectives*, 24(1), 1–21.
- Lewis, D. (1979). Attitudes de dicto and de se. *Philosophical Review*, 88(4), 513–543.
- Lewis, P. J. (2013). The doomsday argument and the simulation argument. *Synthese*, 190(18), 4009–4022.
- Schoenfield, M. (2012). Chilling out on epistemic rationality: A defense of imprecise credences. *Philosophical Studies*, 158(2), 197–219.
- Thomas, T. (2021). Doomsday and objective chance. Global Priorities Institute Working Paper No. 16–2021.
- Tomasik, B. (2016). How the simulation argument dampens future fanaticism. <https://longtermrisk.org/files/how-the-simulation-argument-dampens-future-fanaticism.pdf>
- Weatherson, B. (2003). Are you a sim? *The Philosophical Quarterly*, 53(212), 425–431.
- Weatherson, B. (2005). Should we respond to evil with indifference? *Philosophy and Phenomenological Research*, 70(3), 613–635.
- Williamson, T. (1997). Knowledge as evidence. *Mind*, 106(424), 1–25.

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.