

conflict with the PDE. Here the average utility of an action is simply defined as the sum of utilities of possible worlds where the action is performed divided by their number. To be more precise, let  $b$  be the conjunction of literals true at a world  $w$  in a model  $\mathfrak{M}$ , i.e. the conjunction of the set  $\{\mathfrak{M}, w \models l \mid l \text{ is a literal}\}$ . These are the basic events true in the situation regardless of what or who caused them to be true. We define the overall value of a situation (or world) as  $u(b)$  and write it  $u(\mathfrak{M}, w)$ .

Let  $\mathfrak{M}$  be a model and  $m$  be the number of possible worlds where the action is performed. The average utility of an action is defined as  $(\sum\{u(\mathfrak{M}, w) \mid \mathfrak{M}, w \models a_i\})/m$ . Alternatively, expected utility could be calculated, when probabilities of background conditions are available, or average utility of most plausible worlds, when a plausibility ordering of background conditions is available, but these alternatives are beyond the scope of this paper.

## 6. Dilemma Formalized

We are now in a position, where we can formalize the dilemma presented in [1] and discussed in the beginning of this paper. I will introduce a graphical representation of the causal networks and provide formal details. The arrows between boxes in Figure 1 represent the causal connections between the events. The diamond shaped boxes represent actions. We consider the case where there are two persons to save, H1 and H2. There are three actions in the situation,  $a_1$ , saving H1,  $a_2$ , saving H2, and  $a_3$ , remaining inactive. It is not possible in this situation to save both persons. The consequences of the situations are  $c_1$ , H1 is saved,  $c_2$ , H1 feels discomfort (from being stopped by the robot),  $c_3$ , H2 is saved,  $c_4$ , H2 feels discomfort. We consider just one background condition,  $b_1$ , there are people to be saved. Formally we can specify the dilemma as follows, where we assume the utility of the negation of an event to be the result of multiplying with -1.

$\mathfrak{M} = \langle A = \{a_1, a_2, a_3\}, B = \{b_1\}, C = \{c_1, c_2, c_3, c_4\}, I_1 = \{a_1, c_1\}, I_2 = \{a_2, c_3\}, I_3 = \{a_3\}, N = \{(a_1, c_1), (a_1, c_2), (a_2, c_3), (a_2, c_4), (b_1, c_1), (b_1, c_2), (b_1, c_3), (b_1, c_4)\}, f_1, f_2, f_3, f_4, u(a_1) = u(a_2) = u(a_3) = u(b_1) = 0, u(c_1) = 10, u(c_2) = -4, u(c_3) = 10, u(c_4) = -4, W \rangle$ .

The causal mechanisms  $f_1, f_2, f_3, f_4$  are of level 1 and binary and depend on the background variable and  $a_1, a_2$ ,  $f_1(b_1, a_1), f_2(b_1, a_1), f_3(b_1, a_2), f_4(b_1, a_2)$ . They can be specified as follows.

$$\begin{aligned} c_1 \quad & f_1(1, 1) = 1, f_1(1, 0) = 0, f_1(0, 1) = 1, f_1(0, 0) = 1 \\ c_2 \quad & f_2(1, 1) = 1, f_2(1, 0) = 0, f_2(0, 1) = 0, f_2(0, 0) = 0 \\ c_3 \quad & f_3(1, 1) = 1, f_3(1, 0) = 0, f_3(0, 1) = 1, f_3(0, 0) = 1 \\ c_4 \quad & f_4(1, 1) = 1, f_4(1, 0) = 0, f_4(0, 1) = 0, f_4(0, 0) = 0 \end{aligned}$$

Assuming that there are people to be saved ( $b_1$  is true), we only consider  $W$  to consist of three possible worlds  $w_1, w_2, w_3$ , where  $a_1, a_2, a_3$  are performed respectively. We check the conditions of the PDE for  $a_1$ ,  $a_2$  is similar. At  $w_1$ ,  $cons_1 = \{c_1, c_2\}$ .

1.  $\mathfrak{M}, w_1 \models u(a_1) \geq 0$ .
2. (a)  $\mathfrak{M}, w_1 \models I_{c_1} \rightarrow (u(c_1) \geq 0)$   
(b)  $\mathfrak{M}, w_1 \models I_{c_1} \wedge (u(c_1) > 0)$
3.  $\mathfrak{M}, w_1 \models \neg((c_1 \rightsquigarrow c_2) \wedge ((0 > u(c_1)) \wedge (u(c_2) > 0))) \wedge \neg((c_2 \rightsquigarrow c_1) \wedge ((0 > u(c_2)) \wedge (u(c_1) > 0)))$

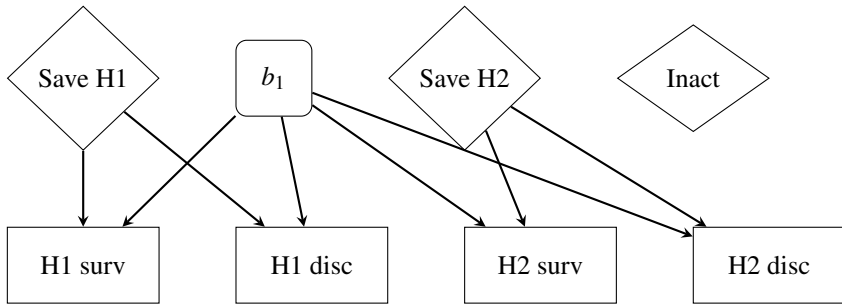


Figure 1. Robot Dilemma

$$4. \mathfrak{M}, w_1 \models u((c_1 \wedge c_2)) > 0$$

Since  $a_1$  meets the conditions of the PDE it is permitted at  $w_1$  and since  $w_1$  is the only world where  $a_1$  is performed it is permitted in the model. The action of remaining inactive,  $a_3$  is also permitted. However as the average utility of performing either  $a_1$  ( $10 \cdot 4 - 10 \cdot 1 = -4$ ) or  $a_2$  ( $10 \cdot 4 - 10 \cdot 1 = -4$ ) is higher than that of  $a_3$  ( $-10 \cdot 10 \cdot 1 = -20$ ), staying inactive, the agent ought to save H1 or H2. Which of  $a_1$  or  $a_2$  to choose can be decided heuristically, as they are both permitted.

## 7. Conclusion

I have shown that the PDE can be formalized and used to resolve certain ethical dilemmas. Even if one does not agree with the principle for philosophical reasons, one might see it as an important effort to bring non-consequentialist aspects of ethical reasoning into the field of robot planning. This is in line with Arkin's idea of the ethical governor, see [4], an ethical safety mechanism that narrows the scope of available actions of social robots and brings us closer to the ideal of an ethical robot. A suitable deontic logic could be used as a meta-reasoning tool, I suggest one in [16]. The implementation e.g. via model checking software, as well as experiments and simulations with social robots are other big research tasks left. So is the challenging task of integrating the PDE with other relevant principles, e.g. from rights-based ethics, with the goal of reaching a reflective equilibrium in each specific case.

## References

- [1] Alan FT Winfield, Christian Blum, and Wenguo Liu. Towards an ethical robot: internal models, consequences and ethical action selection. In M. Mistry, A. Leonardis, M. Witkowski, and C. Melhuish, editors, *Advances in Autonomous Robotics Systems*, pages 85–96. Springer, 2014.
- [2] Selmer Bringsjord and Joshua Taylor. The divine-command approach to robot ethics. In Patrick Lina, Keith Abney, and George A. Bekey, editors, *Robot Ethics: The Ethical and Social Implications of Robotics*, pages 85–108. MIT Press, 2012.
- [3] Wendel Wallach and Colin Allen. *Moral Machines: Teaching Robots Right From Wrong*. Oxford University Press, 2009.
- [4] Ronald Arkin. *Governing Lethal Behavior in Autonomous Robots*. CRC Press, 2009.

- [5] Martin Mose Bentzen. The limits of logic-based inherent safety of social robots. In Diane P. Michelfelder, Byron Newberry, and Qin Zhu, editors, *Philosophy and Engineering : Exploring Boundaries, Expanding Connections*. Springer, forthcoming.
- [6] Philippa Foot. The problem of abortion and the doctrine of double effect. *Oxford Review*, 5:5–15, 1967.
- [7] Judith Jarvis Thomson. The trolley problem. *The Yale Law Journal*, 94:1395–1415, 1985.
- [8] Nick Bostrom. Ethical issues in advanced artificial intelligence. In I. Smit, W. Wallach, and G. Lasker, editors, *Cognitive, Emotive and Ethical Aspects of Decision Making in Humans and in Artificial Intelligence*, volume 2, pages 12–17. Int. Institute of Advanced Studies in Systems Research and Cybernetics, 2003.
- [9] Joseph Mangan. An historical analysis of the principle of double effect. *Theological Studies*, 10:41–61, 1949.
- [10] Warren Quinn. Actions, intentions, and consequences: The doctrine of double effect. *Philosophy and Public Affairs*, 18:334–351, 1989.
- [11] Alison McIntyre. Doctrine of double effect. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Winter 2014 edition, 2014.
- [12] John F. Harty. *Agency and Deontic Logic*. Oxford University Press, 2001.
- [13] Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, 2nd edition, 2009.
- [14] Martin Mose Bentzen. *Stit, Iit, and Deontic Logic for Action Types*. PhD thesis, Section for Philosophy and Science Studies, Roskilde University, 2010.
- [15] Patrick Blackburn, Maarten de Rijke, and Yde Venema. *Modal Logic*. Cambridge University Press, 2001.
- [16] Martin Mose Bentzen. Action type deontic logic. *Journal of Logic, Language and Information*, 23:397–414, 2014.

# What Your Computer Still Can't Know: A Refutation of Bringsjord's Refutation of Searle's Refutation of Bostrom and Floridi

Ahti-Veikko PIETARINEN <sup>a,1</sup>

<sup>a</sup> *Ragnar Nurkse School of Innovation and Governance, Tallinn University of Technology, Estonia*

**Abstract.** I refute Bringsjord's attempted refutation of Searle, who has argued against two recent visions: Bostrom's super-intelligence (post-humanism) and Floridi's info-spheres (information revolution). My refutation derives from the impossibility of Turing machines to compute consequential information not linked with observations of its output. Placing post-humanism and information revolution under a philosophical perspective leads to an identification of an unspoken presupposition in both: universalism of meaning. A philosophical theory of information needs a semiotic theory of signs and representations that take information to be a property of signs that are linked with their interpreting minds.

**Keywords.** philosophy of mind; super-intelligence; post-humanism; info-spheres; information revolution; Searle; philosophy of information; philosophy of signs.

## 1. Introduction

Mental states and their content – conscious or non-conscious – are intentional: they *mean* something. Can they be achieved by computational processes? The following famous objection applies [1]: Computational states are defined entirely in terms of their syntax. That is, computation is formal symbol manipulation: typically electric currents in switching circuits that correspond to 0's and 1's. Computational states cannot be identical with mental states, because mental states also have semantics: the content of thoughts and experiences, their meaning, interpretation and our understanding of that content. That is, syntax is neither equivalent nor sufficient for semantics to arise.

This argument forms the core of Searle's recent [2] refutation of two recent contentions: Bostrom's [3] near-imminent arrival of a potentially malicious super-intelligence, and Floridi's [4] present-day informational revolution and the world as an all-encompassing info-sphere. Both claims appear rather hyperbolic. To refute the former, Searle argues – evoking here the Chinese Room Argument – that computing machines merely manipulate symbols and so cannot be conscious, for to be malicious one would have to be conscious. To establish the implausibility of the world as an info-

---

<sup>1</sup> Chair of Philosophy, Ragnar Nurkse School of Innovation and Governance, Tallinn University of Technology, Estonia. E-mail: ahti-veikko.pietarinen@ttu.ee.

sphere and thus to refute the latter, he applies the property of the observer-relative nature of information.

Bringsjord [5] has hurried to defend both Bostrom and Floridi, by attempting to refute Searle's latest attack on information revolution and machine intelligence evangelists. Bringsjord questions, among others, Searle's assumption that maliciousness presumes consciousness.

I refute Bringsjord's attempt. Hence the validity of Bostrom's and Floridi's assertions remains contested by Searle's arguments. My refutation derives from the impossibility of Turing machines to produce consequential information not linked with observations of its computational output. I then place post-humanism and information revolution under a philosophical perspective, and point out an unspoken presupposition concerning universalism of meaning in both. I propose that a philosophical theory of information benefits from a semiotic theory of signs and an account of representations that take information to be a property of signs that are linked with their interpreting minds.

## **2. Searle's Refutation of Bostrom and Floridi**

Searle presents two arguments: the first is to refute Bostrom's super-intelligence concept and the second to cast serious doubt on Floridi's vision of the world as an info-sphere. Let us deal with each of them in turn.

### *2.1. Super-Intelligence*

According to Bostrom [3], we should be deeply concerned about the possible future arrival of super-intelligent, malicious computing machines, since we might well be targets of their malice.

The emphasis, as I see it, ought to be on "the possible": there is any number of possible anything that may loom large upon us, for the good as much as for the evil. Bostrom states in his book that "if machine brains surpassed human brains in general intelligence, then this new superintelligence could become extremely powerful" – but neither he nor anyone else has presented any plausible accounts even for a remote verity of the antecedent. I do have a further remark to be made on such blind-spots of reasoning that we find in Bostrom's work, in the section that follows my refutation of Bringsjord's criticism of Searle's refutation.

### *2.2. Info-Spheres*

Floridi [4], in turn, has written that the universe in which we live is rapidly populated by vast numbers of information-processing machines, whose level of intelligence, relative to ours, is extremely high. We increasingly understand the universe (including specifically ourselves) through informational conceptualizations. A transformation of all information into a global info-sphere, a massive system of all the data, texts, multimedia, expressions, sounds and memories; all that exists and all that ever were, is no future: it is already happening. We are all part of that global sphere of information. Maximally, Floridi states, the info-sphere is "synonymous with reality".

According to Searle's strategy, Floridi's claim concerning what we mean by information and technology, and hence what the very methods of how the philosophy of ICT is to be conducted, is contestable on similar grounds as Bostrom's.

### 2.3. *Searle's Refutation of Bostrom*

Searle [2] presents the following refutation of Bostrom's claim of the near-future arrival of malicious super-intelligence:

- (1) Computing machines merely manipulate symbols. A machine that manipulates symbols is a syntactic engine. A syntactic engine cannot be conscious.
- (2) A malicious computing machine would by definition be a conscious machine.
- (3) Therefore, no malicious computing machine can exist, let alone arrive on planet Earth.

### 2.4. *Searle's Refutation of Floridi*

In the same review piece, Searle considers Floridi's claim concerning the fourth revolution:

an information revolution so that we all now live in the infosphere ... contains a confusion. The other three revolutions all identify features that are observer independent. ... [But t]he information in question is almost entirely in our attitudes; it is observer relative. [2]

He then presents the following refutation of Floridi's claim about info-spheres:

- (1) Information (unlike the cases that were central to Copernican, Darwinian and Freudian revolutions) is observer-relative.<sup>2</sup>
- (2) Therefore, Floridi's claim about the info-spheres is false.

As such this is an enthymeme, but the missing premises are easily supplied. Searle's notion is that it is only for a conscious agency that information can bear sensible content and meaning. Clearly it is the philosophical concept of information that now distinguishes Searle's view from those of Bostrom's and Floridi's.

## 3. Bringsjord's Refutation of Searle

### 3.1. *A Refutation of Searle's Refutation of Bostrom*

Bringsjord's [4] is an attempt to defend Bostrom and Floridi against Searle's attack. He aims at showing that Searle's argument against Bostrom is false by denying that its second premise would hold ("A malicious computing machine would by definition be a

---

<sup>2</sup>Turing's revolution was not to downplay humans and their centrality in interpreting information, quite the contrary. In fact, Searle, in his response to Floridi's response to Searle's refutation published in the December 18 issue of the *New York Review of Books*, takes Floridi to misattribute Turing the *soi-disant* revolution that Turing himself would not have approved.

conscious machine"). He proposes the following scenario, the actuality of which he takes to be sufficient to show that the premise indeed is false and therefore Searle's argument attacking Bostrom's claim invalid. (I have modified the violent story in some minor details to make it less offensive. This is inessential to the argument):

The year is 2015. An intelligent robot R has just shot and killed an innocent woman in Paris. Before killing the woman, the robot shouted, "I hate non-believing humans!" R then raised its bullet-firing arm and shot the human. R then said, "One less non-believer walking on my streets!" An investigation disclosed that all the relevant internal symbols in R's knowledge-base match perfectly with structures of malice as defined in AI: A formula expressing that R desires to kill the person is there. A formula expressing that R intends to kill the person is there. A formula expressing that R knows how to kill the person is there. The same is found with respect to R's knowledge about who the victim is (a non-believer). And so on. In short, these constitute a definition D' of malice, which says that a robot is malicious if it, as a matter of its internal data, desires to harm innocent people for reasons having nothing to do with preventing further harm etc. The formulation of D' was guided by definitions of malice found in the philosophical literature. (Cf. [4, p. 7])

Bringsjord aims at convincing us that that such robot behaviour is no sci-fi: it is already here, among us. The robots in question have been programmed and constructed in the research projects he himself has founded. And they are, he maintains, what [6] terms *access-conscious*. They are access-conscious by virtue of the formal structure and the reasoning and decision-making capacities that they possess. Bringsjord admits that they are not *phenomenally* conscious, however, which is what Searle asks them to be in order to count as truly malicious or truly informational.

Thus Searle is, according to Bringsjord, wrong about the necessary presence of consciousness for something to count as a malicious piece of hardware.

### 3.2. *Bringsjord's Refutation of Searle's Refutation of Floridi*

Bringsjord's argument against Searle's refutation of Floridi's info-spheres runs along similar considerations. That is, he takes there to exist machines and robots that manipulate information in an observer-independent fashion. For these reasons Bringsjord takes it that Searle's refutation does not go through.

## 4. A Refutation of Bringsjord's Refutation of Searle's Refutation

Bringsjord's refutation is first and foremost an appeal to an empirical observation. But it also has an implicit and much more interesting and interconnected theoretical part.

To begin with the empirical argument, Bringsjord takes a machine's actual existence and its behavioural capacities of performing certain tasks, such as understanding language, to sufficiently demonstrate their capability for autonomous or independent manipulating of information. If so, then information would be observer-independent, that is, independent of the characteristics that we, as we observe those machines and robots in action, would ascribe them to have by our interpretations of what that information means.

But do machines perform tasks such as understand language, in an autonomous and observer-independent fashion? This inevitably takes us to the theoretical and not so well covered areas of the debate. Bringsjord takes the examination of the key attributes of the formal systems that robots are programmed with (or in the very least have been programmed in their research projects) sufficient to guarantee genuine and real understanding.

Those systems are essentially fragments of first-order deontic logic. A question that arises is whether such a logic can support truly autonomous moral reasoning. A problem is that one needs to create and establish a good general ethical code that spells out those fundamental principles on which a machine could rest its decisions on obligations and permissions of its actions whenever it meets with a certain kind of circumstance that calls for deontic resolution. Such a code is hardly universal, as casuist and pragmatistic theories of moral reasoning can be adapted to demonstrate [7].

But what really matters is what the *strategies of reasoning* are. Those strategies would require reading not merely a formal code or a set of principles but a well-defined *semantics*, as well as a grasp of the *pragmatics* of the meaning of actions and behaviour upon which the relevant information is presented to the machine.

My refutation thus boils down to this. Searle has consistently spoken of the *semantic* notion of information as the observer-relative notion. That semantic notion is not something that can adequately be treated either as an *a priori* notion of mathematics or as a casual component of formalized systems of deontic or other modal logics.

Bringsjord, however, can defend only the *syntactic* notion of information, such as whether  $2 + 2 = 4$  indeed holds without the presence of any observers. He explicitly claims that the fact that the truth of elementary arithmetical propositions follows deductively from the Peano axiom systems is “part of the furniture of our universe, even if there be no conscious agents” [4, p. 8].

Now who invented the Peano axioms, one might rhetorically ask here. Who established the main ideas and systems of syllogistics, deductive inference, and such systems of reasoning? These questions alone should suffice to point out where the fundamental mistake lies.

We must also not shun Turing, to whom all parties have made their various appeals. Searle holds that Turing computation and its realizations are observer relative. Bringsjord’s “logico-mathematics” has this “flatly false” by an appeal to what clearly is a Platonic eternity of fundamental theorems of some “logico-mathematics”.

It is beyond dispute, however, as I will now argue, that also Turing machines need an observer for their computations to be interpreted and to make any sense at all. My claim thus differs from Searle’s in that it does not assert that it is a *human* observer or even that it is a *conscious* (in the phenomenal sense) observer that is needed. I only claim that a computing tape or a mathematical proposition not connected to an interpreting mind is inconsequential and ultimately meaningless. I do not claim that such an interpreting mind could not be an artefact, even of a computational kind. What the argument implies is that we do not have such things – that is we do not have them yet, and currently they are not among us.

Whatever the results are that a computational device produces, those results must be observed and interpreted as meaningful. They must be consequential and have practical bearings. This follows from the nature of fundamental truths of mathematics, which are observational, and theoretical computer science hardly makes an exception.



Therefore, Bringsjord's attempted refutation of Searle's argument is ineffective. If Searle's argument holds, and as what the above argument shows it does, then also Floridi's argument is subject to the charge that concerns the observer-relative nature of meaningful information. I will clarify next the main details and some consequences of this argumentative situation.

## 5. Post-humanism in a Philosophical Perspective

Bringsjord seems to err in taking good logic to be equal with a formal system that is programmable and manipulated by a set of generic rules. He even thinks that elementary statements of arithmetic are essentially statements that can be made sensible and understandable within the frameworks of such formal systems.

This theoretical error is at bottom to confuse logical (formal) and non-logical (mathematical) axiom systems. The term "logico-mathematics" betrays him, as do the "mechanized" approaches to his logic and reasoning in AI. What really matters is what interpretations are assigned even to the simplest analytic statements falling from those axiomatic systems. What makes axiomatic systems in mathematics truly useful is not that they are precise rules but that they can be re-interpreted by assigning new meanings to the constants of those axioms. By varying such interpretations mathematicians are able to study different properties of new classes of models thus generated. A non-observer-relative meaning of axioms would simply mean that there could be a non-re-interpretable axiom system. But such a system would have no interesting properties that would emerge from its use and application.

What does this seemingly abstruse point concerning the philosophy of mathematics and logic have to do with those apparently amiable claims often made in lieu of super-intelligence, including radical enhancement or post-humanism?

Everything. To see this, notice that the conceptual issue involved here is two-fold. First, to truly believe in super-intelligence one needs to believe in the notion of information that is observer-independent, that is, information that need not be linked to *any mind* – not only to a human or a conscious mind, that is, but to anything whatever that is mind-like.<sup>3</sup> The latter belief, that is, the presence of information that need not be linked to any mind can, as I argue in the section below, be shown to be implausible by an argument from semiotics.

Second, a closer analysis of those radical claims for singularity and its ilk turns up not only errors in the structure of their arguments (such as missing premises or non-exclusionary either-or arguments) but also, and more interestingly, uncovers an important unspoken presupposition. This presupposition concerns two ways in which our language, as well as conceptual and axiomatic systems, get their meaning.

First, conceptual systems can be thought of as universal: everything I need to know about the meanings in that system can be gleaned from the relationships that obtain in that system. Second, systems can be *re-interpretable*, in which case the meanings are derived from the connections that their interconnected components and expressions have with the universes they discourse about. The distinction between the two can also be termed the *one-world* vs. *many-worlds* hypothesis about meaning (cf. the related van Heijenoort–Hintikka distinction).

---

<sup>3</sup> It is reasonable to request mind-likeliness to have some biological and evolutionary basis, however.

Bringsjord's logic has what is the consequence of the presupposition of the one-world conception of meanings, namely universalism. This is shown by the fact that the deontic logic he appeals to is missing out on its semantic side: it does not tell what the relationships are between the expressions and proofs in one's logic on the one hand and their meanings on the other. It does not tell how actions and interpretations can vary from one circumstance or a model to another. There are no possible-worlds and no situation-theoretic or neighbourhood models, let alone dynamic semantics. The semantics that was proposed uses the utilitarian "deontic stit frames". Even so, such frames are silently encoded into an automated proof-theoretic system with its usual and well-known limitations [8].

As long as such a dynamic, strategic perspective to action is lacking from the operation of internal systems of robots; artefacts that are solely controlled by formal, rule-based deontic reasoning systems do not manipulate information on a meta-theoretical level and in observer-independent fashion that would be required of them by Bringsjord's own arguments.

It is thus this distinction between one and many-worlds approach to meaning that sharply separates Searle from Bringsjord, Bostrom and Floridi. Searle noted how behaviourism and dualism have not disappeared from cognitive sciences. I have added an observation that these residues show up in Bostrom's and Floridi's conceptualizations under the guises of universalism and the one-world presupposition concerning meaning. Bostrom and Floridi both take claims of reality as informational and for that reason are misled to take that vision of reality to be such as to give rise either to autonomous, malicious, uncontrollable, super-intelligent machines, or to the inescapable and closed sphere of information. But such fancies are the results of the universalist presupposition concerning how information gets its meaning in the first place. They are not about various ways in which our assertions relate to the world and to the mind. They are about internal conditions governing the relations of assertions with one another, either within formal systems or within a 'sphere' that already has been declared to be 'informational'.

I refrain from further listing the fallacies that plague the singularity prose. Just to take the simulation argument as another example. Bostrom claims that *either* (i) nearly all human-level civilizations go extinct before becoming post-human, *or* (ii) any post-human civilization is extremely unlikely to run a significant number of simulations of their evolutionary history, *or* (iii) we are almost certainly living in a computer simulation. This is a perfect example of a false trilemma. Our very own world serves as a counterexample. But the burden of proof is on the other side: extraordinary claims need to have extraordinary evidence.

## 6. Information Revolution in a Philosophical Perspective

Floridi invites us to an info-sphere that abounds with ultra-smart technologies. Humans are no longer in its centre, and are likely not in charge of it, either, as the tools and artifacts outperform us in all affairs from routine to more sophisticated ones.

But so did it happen with the tools and instruments or our hominid ancestors: they learned to poke honey from the beehives without getting their arms stung. Seeing the whole universe as computationally informational – phenomenally and relationally, yes – but through and through, with or without human observers, is a tough claim to defend. A computational theory of information, just as the explanation of the operation