

Philosophy of AI

A Structured Overview

Vincent C. Müller

2.1 TOPIC AND METHOD

2.1.1 *Artificial Intelligence*

The term *Artificial Intelligence* became popular after the 1956 “Dartmouth Summer Research Project on Artificial Intelligence,” which stated its aims as follows:

The study is to proceed on the basis of the conjecture that every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it.¹

This is the ambitious research program that human intelligence or cognition can be understood or modeled as rule-based computation over symbolic representation, so these models can be tested by running them on different (artificial) computational hardware. If successful, the computers running those models would display artificial intelligence. Artificial intelligence and cognitive science are two sides of the same coin. This program is usually called *Classical AI*:²

a) AI is a research program to create computer-based agents that have intelligence.

The terms *Strong AI* and *Weak AI* as introduced by John Searle stand in the same tradition. *Strong AI* refers to the idea that: “the appropriately programmed computer really is a mind, in the sense that computers given the right programs can be literally said to understand and have other cognitive states.” *Weak AI* means that AI merely simulates mental states. In this weak sense “the principal value of the computer in the study of the mind is that it gives us a very powerful tool.”³

¹ John McCarthy et al., “A proposal for the Dartmouth Summer Research Project on Artificial Intelligence” (1955), www-formal.stanford.edu/jmc/history/dartmouth/dartmouth.html, accessed October 1, 2006.

² As a sample: Eric Dietrich, “Philosophy of artificial intelligence” *The Encyclopedia of Cognitive Science* 203. The classic historical survey is Margaret A. Boden, *Mind as Machine: A History of Cognitive Science* (Oxford University Press, 2006).

³ John R. Searle, “Minds, brains and programs” (1980) *Behavioral and Brain Sciences*, 3(3): 417–424.

On the other hand, the term “AI” is often used in computer science in a sense that I would like to call *Technical AI*:

- b) AI is a set of computer-science methods for perception, modelling, planning, and action (search, logic programming, probabilistic reasoning, expert systems, optimization, control engineering, neuromorphic engineering, machine learning (ML), etc.).⁴

There is also a minority in AI that calls for the discipline to focus on the ambitions of (a), while maintaining current methodology under (b), usually under the name of *Artificial General Intelligence* (AGI).⁵

This existence of the two traditions (classical and technical) occasionally leads to suggestions that we should not use the term “AI,” because it implies strong claims that stem from the research program (a) but have very little to do with the actual work under (b). Perhaps we should rather talk about “ML” or “decision-support machines,” or just “automation” (as the 1973 Lighthill Report suggested).⁶ In the following we will clarify the notion of “intelligence” and it will emerge that there is a reasonably coherent research program of AI that unifies the two traditions: *The creation of intelligent behavior through computing machines*.

These two traditions now require a footnote: Both were largely developed under the notion of *classical AI*, so what has changed with the move to ML? Machine learning is a traditional computational (connectivist) method in neural networks that does not use representations.⁷ Since ca. 2015, with the advent of massive computing power and massive data for deep neural networks, the performance of ML systems in areas such as translation, text production, speech recognition, games, visual recognition, and autonomous driving has improved dramatically, so that it is superior to humans in some cases. Machine learning is now the standard method in AI. What does this change mean for the future of the discipline? The honest answer is: We do not know yet. Just like any method, ML has its limits, but these limits are less restrictive than was thought for many years because the systems exhibit a non-linear improvement – with more data they may suddenly improve significantly. Its

⁴ Stuart Russell, *Human Compatible: Artificial Intelligence and the Problem of Control* (Viking, 2019); Stuart Russell and Peter Norvig, *Artificial Intelligence: A Modern Approach* (4th ed., Prentice Hall, 2020); Günther Görz, Ute Schmid, and Tanya Braun, *Handbuch der Künstlichen Intelligenz* (5th ed., De Gruyter, 2020); Judea Pearl and Dana Mackenzie, *The Book of Why: The New Science of Cause and Effect* (Basic Books, 2018).

⁵ AGI conferences have been organized since 2008.

⁶ James Lighthill, *Artificial Intelligence: A General Survey* (Science Research Council, 1973).

⁷ Frank Rosenblatt, “The Perceptron: a perceiving and recognizing automaton (Project PARA)” Vol. 85, Issues 460–461 of Report (Cornell Aeronautical Laboratory, 1957); Yann LeCun, Yoshua Bengio, and Geoffrey Hinton, “Deep learning” (2015) *Nature*, 521: 436–444; James Garson and Cameron Buckner, “Connectionism” in Edward N. Zalta (ed.), *Stanford Encyclopedia of Philosophy* (CSLI, Stanford, 2019), <https://plato.stanford.edu/entries/connectionism/>; Cameron J. Buckner, *From Deep Learning to Rational Machines: What the History of Philosophy Can Teach Us about the Future of Artificial Intelligence* (Oxford University Press, 2024).

weaknesses (e.g., overfitting, causal reasoning, reliability, relevance, and black box) may be quite close to those of human rational choice, especially if “predictive processing” is the correct theory of the human mind (Sections 2.4 and 2.6).

2.1.2 *Philosophy of AI and Philosophy*

One way to understand the philosophy of AI is that it mainly deals with three Kantian questions: What is AI? What can AI do? What should AI be? One major part of the philosophy of AI is the *ethics* of AI but we will not discuss this field here, because there is a separate entry on “Ethics of AI” in the present CUP handbook.⁸

Traditionally, the philosophy of AI deals with a few selected points where philosophers have found something to say about AI, for example, about the thesis that cognition is computation, or that computers can have meaningful symbols.⁹ Reviewing these points and the relevant authors (Turing, Wiener, Dreyfus, Dennett, Searle, ...) would result in a fragmented discussion that never achieves a picture of the overall project. It would be like writing an old-style human history through a few “heroes.” Also, in this perspective, the philosophy of AI is separated from its cousin, the philosophy of cognitive science, which in turn is closely connected to the philosophy of mind.¹⁰

In this chapter we use a different approach: We look at *components of an intelligent system*, as they present themselves in philosophy, cognitive science, and AI. One way to consider such components is that there are relatively simple animals that can do relatively simple things, and then we can move “up” to more complicated animals that can do those simple things, and more. As a schematic example, a *fly* will continue to bump into the glass many times to get to the light; a *cobra* will understand that there is an obstacle here and try to avoid it; a *cat* might remember that there was an obstacle there the last time and take another path right away; a *chimpanzee* might realize that the glass can be broken with a stone; a *human* might find the key and unlock the glass door ... or else take the window to get out.

⁸ See Chapter 3 of this book. See also: Vincent C. Müller, “Ethics of artificial intelligence and robotics” in Edward N. Zalta (ed.), *Stanford Encyclopedia of Philosophy*, (CSLI, Stanford University, 2020), <https://plato.stanford.edu/entries/ethics-ai/>; Vincent C. Müller, *Can Machines Think? Fundamental Problems of Artificial Intelligence* (Oxford University Press, forthcoming).

⁹ There are very few surveys and no recent ones. See Jack B. Copeland, *Artificial Intelligence: A Philosophical Introduction* (Blackwell, 1993); Dietrich Matt Carter, *Minds and Computers: An Introduction to the Philosophy of Artificial Intelligence* (Edinburgh University Press, 2007); Luciano Floridi (ed.) *The Blackwell Guide to the Philosophy of Computing and Information* (Blackwell, 2003); Luciano Floridi, *The Philosophy of Information* (Oxford University Press, 2011). Some of what philosophers had to say can be seen as undermining the project of AI, compare Eric Dietrich et al., *Great Philosophical Objections to Artificial Intelligence: The History and Legacy of the AI Wars* (Bloomsbury Academic, 2021).

¹⁰ Eric Margolis, Richard Samuels, and Stephen Stich (eds.), *The Oxford Handbook of Philosophy of Cognitive Science* (Oxford University Press, 2012).

To engage in the philosophy of AI properly, we will thus need a wide range of philosophy: philosophy of mind, epistemology, language, value, culture, society, ...

Furthermore, in our approach, the philosophy of AI is not just “applied philosophy”; it is not that we have a solution ready in the philosopher’s toolbox and “apply” it to solve issues in AI. The philosophical understanding itself *changes* when looking at the case of AI: It becomes less anthropocentric, less focused on our own human case. A deeper look at concepts must be normatively guided by the *function* these concepts serve, and that function can be understood better when we consider both the natural cases *and* the case of actual and possible AI. This chapter is thus also a “proof of concept” for doing philosophy through the conceptual analysis of AI: I call this *AI philosophy*.

I thus propose to turn the question from its head onto its feet, as Marx would have said: If we want to understand AI, we have to understand ourselves; and if we want to understand ourselves, we have to understand AI!

2.2 INTELLIGENCE

2.2.1 *The Turing Test*

“I propose to consider the question ‘Can Machines Think?’” Alan Turing wrote at the outset of his paper in the leading philosophical journal *Mind*.¹¹ This was 1950, Turing was one of the founding fathers of computers, and many readers of the paper would not even have heard of such machines, since there were only half a dozen universal computers in the world (Z3, Z4, ENIAC, SSEM, Harvard Mark III, and Manchester Mark I).¹² Turing moves swiftly to declare that searching for a definition of “thinking” would be futile and proposes to replace his initial question by the question whether a machine could successfully play an “imitation game.” This game has come to be known as the “Turing Test”: A human interrogator is connected to another human and a machine via “teleprinting,” and if the interrogator cannot tell the machine from the human by holding a conversation, then we shall say the machine is “thinking.” At the end of the paper he returns to the issue of whether machines can think and says: “I believe that at the end of the century the use of words and general educated opinion will have altered so much that one will be able to speak of machines thinking without expecting to be contradicted.”¹³ So, Turing proposes to replace our everyday term of “thinking” by an operationally defined term, a term for which we can test with some procedure that has a measurable outcome.

Turing’s proposal to replace the definition of thinking by an operational definition that relies exclusively on behavior fits with the intellectual climate of the time

¹¹ Alan Turing, “Computing machinery and intelligence” LIX *Mind* 433.

¹² Anonymous, “Digital computing newsletter” (1950) 2 Office of Naval Research, Mathematical Sciences Division, Washington DC 1.

¹³ Turing 442.

where behaviorism became a dominant force: In psychology, behaviorism is a *methodological* proposal that psychology should become a proper scientific discipline by relying on testable observation and experiment, rather than on subjective introspection. Given that the mind of others is a “black box,” psychology should become the science of stimulus and behavioral response, of an input–output relation. Early analytic philosophy led to *reductionist behaviorism*; so if the meaning of a term is its “verification conditions,” then a mental term such as “pain” just *means* the person is disposed to behaving a certain way.

Is the Turing Test via observable behavior a useful definition of intelligence? Can it “replace” our talk of intelligence? It is clear that there will be intelligent beings that will not pass this test, for example, humans or animals that cannot type. So I think it is fair to say that Turing very likely only intended the passing of the test as being sufficient for having intelligence and not as necessary. So, if a system passes that test, does it have to be intelligent? This depends on whether you think intelligence is just intelligent behavior, or whether you think for the attribution of intelligence we also need to look at internal structure.

2.2.2 What Is Intelligence?

Intuitively, intelligence is an ability that underlies intelligent action. Which action is intelligent depends on the goals that are pursued, and on the success in achieving them – think of the animal cases mentioned earlier. Success will depend not only on the agent but also on the conditions in which it operates, so a system with fewer options how to achieve a goal (e.g., find food) is less intelligent. In this vein, a classical definition is: “Intelligence measures an agent’s ability to achieve goals in a wide range of environments.”¹⁴ Here intelligence is the *ability to flexibly pursue goals*, where flexibility is explained with the help of different environments. This notion of intelligence from AI is an *instrumental* and normative notion of intelligence, in the tradition of classical decision theory, which says that a rational agent should always try to maximize expected utility (see Section 2.6).¹⁵

If AI philosophy understands intelligence as relative to an environment, then to achieve more intelligence, one can change the agent or change the environment. Humans have done both on a huge scale through what is known as “culture”: Not only have we generated a sophisticated learning system for humans (to change the agent), we have also physically shaped the world such that we can pursue our goals in it; for example, to travel, we have generated roads, cars with steering wheels,

¹⁴ Shane Legg and Marcus Hutter, “Universal intelligence: A definition of machine intelligence” (2007) *Minds and Machines*, 17(4): 391–444, 402.

¹⁵ See, for example, Herbert A. Simon, “A behavioral model of rational choice” (1955) *Quarterly Journal of Economics*, 69(1): 99–118; Johanna Thoma, “Decision theory” in Richard Pettigrew and Jonathan Weisberg (eds), *The Open Handbook of Formal Epistemology* (PhilPapers, 2019), see also the neo-behaviorist proposal in Dimitri Coelho Mollo, “Intelligent Behaviour” [2022] *Erkenntnis*.

maps, road signs, digital route planning, and AI systems. We now do the same for AI systems; both the learning system, and the change of the environment (cars with computer interfaces, GPS, etc.). By changing the environment, we will also change our cognition and our lives – perhaps in ways that turn out to be to our detriment.

In Sections 2.4–2.9, we will look at the main components of an intelligent system; but before that we discuss the mechanism used in AI: computation.

2.3 COMPUTATION

2.3.1 *The Notion of Computation*

The machines on which AI systems run are “computers,” so it will be important for our task to find out what a computer is and what it can do, in principle. A related question is whether human intelligence is wholly or partially due to computation – if it is wholly due to computation, as classical AI had assumed, then it appears possible to recreate this computation on an artificial computing device.

In order to understand what a computer is, it is useful to remind ourselves of the history of computing machines – I say “machines” because before ca. 1945, the word “computer” was a term for a human who has a certain profession, for someone who does computations. These computations, for example, the multiplication of two large numbers, are done through a mechanical step-by-step procedure that will lead to a result once carried out completely. Such procedures are called “algorithms.” In 1936, in response to Gödel’s challenge of the “Entscheidungsproblem,” Alan Turing suggested that the notion of “computing something” could be explained by “what a certain type of machine can do” (just like he proposed to operationalize the notion of intelligence in the “Turing Test”). Turing sketched what such a machine would look like, with an infinitely long tape for memory, a head that can read from and write symbols to that tape. These states on the tape are always specific discrete states, such that each state is of a type from a finite list (symbols, numbers,...), so for example it either is the letter “V” or the letter “C,” not a bit of each. In other words, the machine is “digital” (not analog).¹⁶ Then there is one crucial addition: In the “universal” version of the machine, one can *change* what the computer does through further input. In other words, the machine is *programmable* to perform a certain algorithm, and it stores that program in its memory.¹⁷ Such a computer is a universal

¹⁶ Nicholas Negroponte, *Being digital* (Vintage, 1995); see also John Haugeland, *Artificial intelligence: The very idea* (MIT Press, 1985) 57; Vincent C. Müller, “What is a digital state?” in Mark J. Bishop and Yasemin J. Erden (eds), *The Scandal of Computation – What Is Computation? – AISB Convention 2013* (AISB, 2013), www.aishb.org.uk/asibpublications/convention-proceedings.

¹⁷ Alan Turing, “On computable numbers, with an application to the Entscheidungsproblem” 2 Proceedings of the London Mathematical Society 230 Kurt Gödel, “Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme I” 38 Monatshefte für Mathematik und Physik 173. The original program outlined in David Hilbert, “Mathematische Probleme” [Springer] Nachrichten von der Königlichen Gesellschaft der Wissenschaften zu Göttingen, Math-Phys Klasse

computer, that is, it can compute any algorithm. It should be mentioned that wider notions of computation have been suggested, for example, analog computing and hypercomputing.¹⁸

There is also the question whether computation is a real property of physical systems, or whether it is rather a useful way of describing these. Searle has said: “The electrical state transitions are intrinsic to the machine, but the computation is in the eye of the beholder.”¹⁹ If we take an anti-realist account of computation, then the situation changes radically.

The exact same computation can be performed on different physical computers, and it can have a different semantics. There are thus three levels of description that are particularly relevant for a given computer: (a) The *physical level* of the actual “realization” of the computer, (b) the *syntactic level* of the algorithm computed, and (c) the *symbolic level* of content, of what is computed.

Physically, a computing machine can be built out of anything and use any kind of property of the physical world (cogs and wheels, relays, DNA, quantum states, etc.). This can be seen as using a physical system to encode a formal system.²⁰ Actually, all universal computers have been made with large sets of switches. A switch has two states (open/closed), so the resulting computing machines work on two states (on/off, 0/1), they are *binary* – this is a design decision. Binary switches can easily be combined to form “logic gates” that operate on input in the form of the logical connectives in Boolean logic (which is also two-valued): NOT, AND, OR, and so on. If such switches are in a state that can be *syntactically* understood as 1010110, then *semantically*, this could (on current ASCII/ANSI conventions) represent the letter “V,” the number “86,” a shade of light gray, a shade of green, and so on.

2.3.2 Computationalism

As we have seen, the notion that computation is the cause of intelligence in natural systems, for example, humans, and can be used to model and reproduce this intelligence is a basic assumption of classical AI. This view is often coupled with (and

253. See, for example, Jack B. Copeland, Carl J. Posy, and Oron Shagrir, *Computability: Turing, Gödel, Church, and Beyond* (MIT Press, 2013).

¹⁸ Hava T. Siegelmann, “Computation beyond the Turing limit” *Science* 545; *Neural Networks and Analog Computation: Beyond the Turing Limit* (Birkhäuser, 1997); Oron Shagrir, *The Nature of Physical Computation* (Oxford University Press, 2022); Gualtiero Piccinini, “Computation in physical systems” (2010) *Stanford Encyclopedia of Philosophy*, <https://plato.stanford.edu/entries/computation-physicalsystems/>.

¹⁹ John R. Searle, *Mind: A Brief Introduction* (Oxford University Press, 2004); Gordana Dodig-Crnkovic and Vincent C. Müller, “A dialogue concerning two world systems: Info-computational vs. mechanistic” in Gordana Dodig-Crnkovic and Mark Burgin (eds), *Information and Computation: Essays on Scientific and Philosophical Understanding of Foundations of Information and Computation* (World Scientific, 2011), <https://worldscientific.com/worldscibooks/10.1142/7637#t=aboutBook>.

²⁰ Clare Horsman et al., “When does a physical system compute?” (2014) *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Science*, 470(2169): 1–25.

motivated by) the view that human mental states are functional states and that these functional states are that of a computer: “machine functionalism.” This thesis is often assumed as a matter of course in the cognitive sciences and neuroscience, but it is also the subject of significant criticism in recent decades.²¹ The main sources for this view are an enthusiasm for the universal technology of digital computation, and early neuroscientific evidence indicated that human neurons (in the brain and body) are also somewhat binary, that is, either they send a signal to other neurons, they “fire,” or they don’t. Some authors defend the *Physical Symbol System Hypothesis*, which is computationalism, plus the contention that only computers can be intelligent.²²

2.4 PERCEPTION AND ACTION

2.4.1 *Passive Perception*

You may be surprised to find that the heading of this chapter combines perception and action in one. We can learn from AI and cognitive science that the main function of perception is to allow action; indeed that perception *is* a kind of action. The traditional understanding of perception in philosophy is *passive* perception, watching ourselves watching the world in what Dan Dennett has called the *Cartesian Theatre*: It is as though I had a little human sitting inside my head, listening to the outside world through our ears, and watching the outside world through our eyes.²³ That notion is absurd, particularly because it would require there to be yet another little human sitting in the head of that little human. And yet, a good deal of the discussion of human perception in the philosophical literature really does treat perception as though it were something that happens to me when inside.

For example, there is the 2D–3D problem in vision, the problem of how I can generate the visual experience of a 3D world through a 2D sensing system (the retina, a 2D sheet that covers our eyeballs from the inside). There must be a way of processing the visual information in the retina, the optical nerve and the optical processing centers of the brain that generates this 3D experience. Not really.²⁴

²¹ Marcin Miłkowski, “Objections to computationalism: A survey” *Roczniki Filozoficzne*, LXVI: 1; Shimon Edelman, *Computing the Mind: How the Mind Really Works* (Oxford University Press, 2008), for the discussion Stevan Harnad, “The symbol grounding problem” *Physica D*, 42: 335; Matthias Scheutz (ed) *Computationalism: New Directions* (Cambridge University Press, 2002); Oron Shagrir, “Two dogmas of computationalism” *Minds and Machines*, 7: 321; Francisco J. Varela, Evan Thompson, and Eleanor Rosch, *The Embodied Mind: Cognitive Science and Human Experience* (MIT Press, 1991).

²² Allen Newell and Herbert A. Simon, “Computer science as empirical inquiry: Symbols and search” *Communications of the ACM*, 19(3): 113–126, 116; cf. Boden 1419ff.

²³ Dennett D. C., *Consciousness Explained* (Little, Brown & Co., 1991), 107.

²⁴ For an introduction to vision, see Kevin J. O’Regan, *Why Red Doesn’t Sound Like a Bell: Understanding the Feel of Consciousness* (Oxford University Press, 2011), chapters 1–5.

2.4.2 Active Perception

Actually, the 3D impression is generated by an interaction between me and the world (in the case of vision it involves movement of my eyes and my body). It is better to think of perception along with the lines of the sense of touch: Touching is something that I *do*, so that I can find out the softness of an object, the texture of its surface, its temperature, its weight, its flexibility, and so on. I do this by acting and then perceiving the change of sensory input. This is called a perception-action-loop: I do something, that changes the world, and that changes the perception that I have.

It will be useful to stress that this occurs with perception of my own body as well. I only know that I have a hand because my visual sensation of the hand, the proprioception, and the sense of touch are in agreement. When that is not the case it is fairly easy to make me feel that a rubber hand is my own hand – this is known as the “rubber hand illusion.” Also, if a prosthetic hand is suitably connected to the nervous system of a human, then the perception-action-loop can be closed again, and the human will feel this as their own hand.

2.4.3 Predictive Processing and Embodiment

This view of perception has recently led to a theory of the “predictive brain”: What the brain does is not to passively wait for input, but it is *always on* to actively participate in the action-perception-loop. It generates *predictions* what the sensory input will be, given my actions, and then it matches the predictions with the actual sensory input. The difference between the two is something that we try to minimize, which is called the “free energy principle.”²⁵

In this tradition, the perception of a natural agent or AI system is something that is intimately connected to the physical interaction of the body of the agent with the environment; perception is thus a component of embodied cognition. A useful slogan in this context is “4E cognition,” which says that cognition is *embodied*; it is *embedded* in an environment with other agents; it is *enactive* rather than passive; and it is *extended*, that is, not just inside the head.²⁶ One aspect that is closely connected to 4E cognition is the question whether cognition in humans is

²⁵ Andy Clark, “Whatever next? Predictive brains, situated agents, and the future of cognitive science.” *Behavioral and Brain Sciences*, 36: 181; Andy Clark, *Surfing Uncertainty: Prediction, Action, and the Embodied Mind* (Oxford University Press, 2016); Karl J. Friston, “The free-energy principle: A unified brain theory?” *Nature Reviews Neuroscience*, 11: 127.

²⁶ Andy Clark, *Natural Born Cyborgs: Minds, Technologies, and the Future of Human Intelligence* (Oxford University Press, 2003); Andy Clark and David J. Chalmers, “The extended mind” *Analysis*, 58: 7; Albert Newen, Shaun Gallagher, and Leon De Bruin, “4E Cognition: Historical roots, key concepts, and central issues” in Albert Newen, Leon De Bruin, and Shaun Gallagher (eds), *The Oxford Handbook of 4E Cognition* (Oxford Academic, Oxford University Press, 2018), pp. 3–16, <https://doi.org/10.1093/oxfordhb/9780198735410.013.1>, accessed June 29, 2023.

fundamentally representational, and whether cognition in AI has to be representational (see Section 2.5).

Embodied cognition is sometimes presented as an empirical thesis about actual cognition (especially in humans) or as a thesis on the suitable design of AI systems, and sometimes as an analysis of what cognition is and has to be. In the latter understanding, non-embodied AI would necessarily miss certain features of cognition.²⁷

2.5 MEANING AND REPRESENTATION

2.5.1 *The Chinese Room Argument*

As we saw earlier, classical AI was founded on the assumption that the appropriately programmed computer really *is* a mind – this is what John Searle called *strong AI*. In his famous paper “Minds, Brains and Programs,” Searle presented a thought experiment of the “Chinese Room.”²⁸ The Chinese Room is a computer, constructed as follows: There is a closed room in which John Searle sits and has a large book that provides him with a computer program, with algorithms, on how to process the input and provide output. Unknown to him, the input that he gets is Chinese writing, and the output that he provides are sensible answers or comments about that linguistic input. This output, so the assumption, is indistinguishable from the output of a competent Chinese speaker. And yet Searle in the room understands no Chinese and will learn no Chinese from the input that he gets. Therefore, Searle concludes, *computation is not sufficient for understanding*. There can be no strong AI.

In the course of his discussion of the Chinese room argument, Searle looks at several replies: The *systems reply* accepts that Searle has shown that no amount of simple manipulation of the person in the room will enable that person to understand Chinese, but objects that perhaps symbol manipulation will enable *the wider system*, of which the person is a component, to understand Chinese. So perhaps there is a part-whole fallacy here? This reply raises the question, why one might think that the whole system has properties that the algorithmic processor does not have.

One way to answer this challenge, and change the system, is the *robot reply*, which grants that the whole system, as described, will not understand Chinese because it is missing something that Chinese speakers have, namely a causal connection between the words and the world. So, we would need to add sensors and actuators to this computer, that would take care of the necessary causal connection. Searle responds to this suggestion by saying input from sensors would be “just more

²⁷ Hubert L. Dreyfus, *What Computers Still Can't Do: A Critique of Artificial Reason* (2nd ed, MIT Press, 1992, 1972); Rolf Pfeifer and Josh Bongard, *How the Body Shapes the Way We Think: A New View of Intelligence* (MIT Press, 2007).

²⁸ Searle, “Minds, brains and programs.”

Chinese” to Searle in the room; it would not provide any further understanding, in fact Searle would have no idea that the input is from a sensor.²⁹

2.5.2 Reconstruction

I think it is best to view the core of the Chinese room argument as an extension of Searle’s remark:

No one would suppose that we could produce milk and sugar by running a computer simulation of the formal sequences in lactation and photosynthesis, but where the mind is concerned many people are willing to believe in such a miracle.³⁰

Accordingly, the argument that remains can be reconstructed as:

- a. If a system does only syntactical manipulation, it will not acquire meaning.
- b. A computer does only syntactical manipulation.
-
- c. A computer will not acquire meaning.

In Searle’s terminology, a computer has *only syntax* and *no semantics*; the symbols in a computer lack the intentionality (directedness) that human language use has. He summarizes his position at the end of the paper:

“Could a machine think?” The answer is, obviously, yes. We are precisely such machines. [...] But could something think, understand, and so on solely in virtue of being a computer with the right sort of program? [...] the answer is no.³¹

2.5.3 Computing, Syntax, and Causal Powers

Reconstructing the argument in this way, the question is whether the premises are true. Several people have argued that premise 2 is false, because one can only understand what a computer does as responding to the program as meaningful.³² I happen to think that this is a mistake, the computer does not *follow* these rules, it is just constructed in such a way that it *acts according to* these rules, if its states are suitably

²⁹ David Cole, “The Chinese room argument” (2020), <http://plato.stanford.edu/entries/chinese-room/>; John Preston and Mark Bishop (eds), *Views into the Chinese Room: New Essays on Searle and Artificial Intelligence* (Oxford University Press, 2002).

³⁰ Searle, ‘Minds, brains and programs’ 424.

³¹ Ibid.

³² John McCarthy, “John Searle’s Chinese room argument” (2007), www-formal.stanford.edu/jmc/chinese.html, accessed June 10, 2007; John Haugeland, “Syntax, semantics, physics” in John Preston and Mark Bishop (eds), *Views into the Chinese Room: New Essays on Searle and Artificial Intelligence* (Oxford University Press, 2002) 385; Margaret A. Boden, *Computer Models of the Mind: Computational Approaches in Theoretical Psychology* (Cambridge University Press, 1988) 97.

interpreted by an observer.³³ Having said that, any actual computer, any physical realization of an abstract algorithm processor, *does* have causal powers, it does more than syntactic manipulation. For example, it may be able to turn the lights on or off.

The Chinese room argument has moved the attention in the philosophy of language away from convention and logic toward the conditions for a speaker to mean what they say (speakers' meaning), or to mean anything at all (have intentionality); in particular, it left us with the discussion on the role of *representation* in cognition, and the role of computation over representations.³⁴

2.6 RATIONAL CHOICE

2.6.1 Normative Decision Theory: MEU

A rational agent will perceive the environment, find out which options for action exist, and then take the best decision. This is what decision theory is about. It is a normative theory on how a rational agent *should* act, given the knowledge they have – not a descriptive theory of how rational agents *will* actually act.

So how should a rational agent decide which is the best possible action? They evaluate the possible outcomes of each choice and then select the one that is best, meaning the one that has the highest subjective utility, that is, utility as seen by the particular agent. It should be noted that rational choice in this sense is not necessarily egoistic, it could well be that the agent puts a high utility on the happiness of someone else, and thus rationally chooses a course of action that maximizes overall utility through the happiness of that other person. In actual situations, the agent typically does not know what the outcomes of particular choices will be, so they act under uncertainty. To overcome this problem, the rational agent selects the action with *maximum expected utility* (MEU), where the value of a choice equals the utility of the outcome multiplied by the probability of that outcome occurring. This thought can be explained with the expected utility of certain gambles or lotteries. In more complicated decision cases the rationality of a certain choice depends on subsequent choices of *other agents*. These kinds of cases are often described with the help of “games” played with other agents. In such games it is often a successful strategy to cooperate with other agents in order to maximize subjective utility.

In artificial intelligence it is common to perceive of AI agents as rational agents in the sense described. For example, Stuart Russell says: “In short, a rational agent acts so as to maximise expected utility. It’s hard to over-state the importance of this conclusion. In many ways, artificial intelligence has been mainly about working out the details of how to build rational machines.”³⁵

³³ Ludwig Wittgenstein, “Philosophische Untersuchungen” in *Schriften I* (Suhrkamp, 1980, 1953) §82.

³⁴ John R. Searle, “Intentionality and its place in nature” in *Consciousness and Language* (Cambridge University Press, 2002, 1984); Searle, *Mind: A brief introduction*.

³⁵ Russell 23.

2.6.2 Resources and Rational Agency

It is not the case that a rational agent *will* always choose the perfect option. The main reason is that such an agent must deal with the fact that their resources are limited, in particular, data storage and time (most choices are time-critical). The question is thus not only what the best choice is, but how many resources I should spend on optimizing my choice; when should I stop optimizing and start acting? This phenomenon is called *bounded rationality*, *bounded optimality*, and in cognitive science, it calls for *resource rational* analysis.³⁶ Furthermore, there is no set of discrete options from which to choose, and a rational agent needs to reflect on the goals to pursue (see Section 2.9).

The point that agents (natural or artificial) will have to deal with limited resources when making choices, has tremendous importance for the understanding of cognition. It is often not fully appreciated in philosophy – even the literature about the limits of rational choice seems to think that there is something “wrong” with using heuristics that are biased, being “nudged” by the environment, or using the environment for “extended” or “situated” cognition.³⁷ But it would be irrational to aim for perfect cognitive procedures, not to mention for cognitive procedures that would not be influenced by the environment.

2.6.3 The Frame Problem(s)

The original frame problem for classical AI was how to *update a belief system* after an action, without stating all the things that have *not* changed; this requires a logic where conclusions can change if a premise is added – a non-monotonic logic.³⁸ Beyond this more technical problem, there is a philosophical problem of updating beliefs after action, popularized by Dennett, which asks how to find out what is relevant, how wide the frame should be cast for *relevance*.

³⁶ Simon 99. Gregory Wheeler, “Bounded rationality” in Edward N. Zalta (ed), *The Stanford Encyclopedia of Philosophy*, vol Fall 2020 Edition (CSLI, 2020), <https://plato.stanford.edu/archives/fall2020/entries/bounded-rationality/>; Stuart Russell, “Rationality and intelligence: A brief update” in Vincent C. Müller (ed), *Fundamental Issues of Artificial Intelligence* (Springer, 2016) 16ff; Falk Lieder and Thomas L. Griffiths, “Resource-rational analysis: Understanding human cognition as the optimal use of limited computational resources” (Cambridge University Press) 43; *Behavioral and Brain Sciences*, e1.

³⁷ Daniel Kahneman and Amos Tversky, “Prospect theory: An analysis of decision under risk” *Econometrica*, 47: 263; Daniel Kahnemann, *Thinking Fast and Slow* (Macmillan, 2011); Richard H. Thaler and Cass Sunstein, *Nudge: Improving Decisions about Health, Wealth and Happiness* (Penguin, 2008) vs. David Kirsh, “Problem solving and situated cognition” in P. Robbins and M. Aydede (eds), *The Cambridge Handbook of Situated Cognition* (Cambridge University Press, 2009).

³⁸ Murray Shanahan, “The frame problem” in Edward N. Zalta (ed), *Stanford Encyclopedia of Philosophy*, vol Spring 2016 edition (CSLI, Stanford University, 2016), <https://plato.stanford.edu/archives/spr2016/entries/frame-problem/>.

As Shanahan says “relevance is holistic, open-ended, and context-sensitive” but logical inference is not.³⁹

There is a very general version of the frame problem, expressed by Jerry Fodor, who says, the frame problem really is: “Hamlet’s problem: when to stop thinking.” He continues by saying that “modular cognitive processing is *ipso facto* irrational [...] by attending to less than all the evidence that is relevant and available.”⁴⁰ Fodor sets the challenge that in order to perform a logical inference, especially an abduction, one needs to have decided what is relevant. However, he seems to underestimate that one cannot attend to *all* that is relevant and available (rationality is bounded). It is currently unclear whether the frame problem can be formulated without dubious assumptions about rationality. Similar concerns apply to the claims that Gödel has shown deep limitations of AI systems.⁴¹ Overall, there may be more to intelligence than instrumental rationality.

2.6.4 Creativity

Choices that involve *creativity* are often invoked as something special, not merely mechanical, and thus inaccessible to a mere machine. The notion of “creation” has significant impact in our societal practice particularly when that creation is protected by intellectual property rights – and AI systems *have* created or cocreated music, painting, and text. It is not clear that there is a notion of creativity that would provide an argument against machine creativity. Such a notion would have to combine two aspects that seem to be in tension: On the one hand, creativity seems to imply causation that includes acquiring knowledge and techniques (think of J. S. Bach composing a new cantata), on the other hand, creativity is supposed to be a non-caused, non-predictable, spark of insight. It appears unclear whether such a notion of creativity can, or indeed should, be formulated.⁴² Perhaps a plausible account is that creativity involves moving between different spaces of relevance, as in the frame problem.

³⁹ Daniel C. Dennett, “Cognitive wheels: The frame problem of AI” in Christopher Hookway (ed), *Minds, Machines, and Evolution: Philosophical Studies* (Cambridge University Press, 1984).

⁴⁰ Jerry A. Fodor, “Modules, frames, fridgeons, sleeping dogs, and the music of the spheres” in J. L. Garfield (ed), *Modularity in Knowledge Representation and Natural-Language Understanding* (The MIT Press, 1987) 140f; Dan Sperber and Deirdre Wilson, “Fodor’s frame problem and relevance theory” *Behavioral and Brain Sciences*, 19: 530.

⁴¹ J. R. Lucas, “Minds, machines and Gödel: A retrospect” in Peter J. R. Millican and Andy Clark (eds), *Machines and Thought* (Oxford University Press, 1996); Peter Koellner, “On the question of whether the mind can be mechanized, I: From Gödel to Penrose” *Journal of Philosophy*, 115: 337; Peter Koellner, “On the question of whether the mind can be mechanized, II: Penrose’s new argument” *Journal of Philosophy*, 115: 453.

⁴² Margaret A. Boden, “Creativity and artificial intelligence: A contradiction in terms?” in Elliot Samuel Paul and Scott Barry Kaufman (eds), *The Philosophy of Creativity: New Essays* (Oxford University Press, 2014), 224–244, <https://philpapers.org/archive/PAUTPO-3.pdf>; Simon Colton and Geraint A. Wiggins, *Computational Creativity: The Final Frontier?* (Montpelier, 2012); Martha Halina, “Insightful artificial intelligence” *Mind and Language*, 36: 315.

2.7 FREE WILL AND CREATIVITY

2.7.1 *Determinism, Compatibilism*

The problem that usually goes under the heading of “free will” is how physical beings like humans or AI systems can have something like free will. The traditional division for possible positions in the space of free will can be put in terms of a decision tree. The first choice is whether *determinism* is true, that is, the thesis that all events are caused. The second choice is whether *incompatibilism* is true, that is, the thesis that if determinism is true, then there is no free will.

The position known as *hard determinism* says that determinism is indeed true, and if determinism is true then there is no such thing as free will – this is the conclusion that most of its opponents try to avoid. The position known as *libertarianism* (not the political view) agrees that incompatibilism is true, but adds that determinism is not, so we are free. The position known as *compatibilism* says that determinism and free will are compatible and thus it may well be that determinism is true *and* humans have free will (and it usually adds that this is actually the case).

This results in a little matrix of positions:

		Incompatibilism	Compatibilism
Determinism	Hard Determinism	Optimistic/Pessimistic Compatibilism	
Non-Determinism	Libertarianism	[Not a popular option]	

2.7.2 *Compatibilism and Responsibility in AI*

In a first approximation, when I say I did something freely, it means that it was *up to me* that I was *in control*. That notion of control can be cashed out by saying I could have done otherwise than I did, specifically I could have done otherwise if I had *decided* otherwise. To this we could add that I would have decided otherwise if I had had other *preferences* or *knowledge* (e.g., I would not have eaten those meatballs if I had a preference against eating pork, and if I had known that they contain pork). Such a notion of freedom thus involves an *epistemic condition* and a *control condition*.

So, I act freely if I do as I choose according to the preferences that I have (my subjective utility). But why do I have these preferences? As Aristotle already knew, they are not under my voluntary control, I could not just *decide* to have other preferences and then have them. However, as Harry Frankfurt has pointed out, I can have *second-order* preferences or desires, that is, I can prefer to have other preferences than the ones I actually have (I could want not to have a preference for those meatballs, for example). The notion that I can overrule my preferences with rational thought is what Frankfurt calls the *will*, and it is his condition for being a person.

In a first approximation one can thus say, *to act freely is to act as I choose, to choose as I will, and to will as I rationally decide to prefer*.⁴³

The upshot of this debate is that the function of a notion of free will for agency in AI or humans is to allow personal *responsibility*, not to determine *causation*. The real question is: What are the conditions such that an agent is *responsible* for their actions and *deserves* being praised or blamed for them. This is independent of the freedom from causal determination; that kind of freedom we do not get, and we do not need.⁴⁴

There is a further debate between “optimists” and “pessimists” whether humans actually do fulfil those conditions (in particular whether they can truly cause their preferences) and can thus properly be said to be responsible for their actions and *deserve* praise or blame – and accordingly whether reward or punishment should have mainly forward-looking reasons.⁴⁵ In the AI case, an absence of responsibility has relevance for their status as moral agents, for the existence of “responsibility gaps,” and for what kinds of decisions we should leave to systems that cannot be held responsible.⁴⁶

2.8 CONSCIOUSNESS

2.8.1 Awareness and Phenomenal Consciousness

In a first approximation, it is useful to distinguish two types of consciousness: *Awareness* and *phenomenal consciousness*. Awareness is the notion that a system has cognitive states on a base level (e.g., it senses heat) and on a meta level, it has states where it is aware of the states on the object level. This awareness, or access, involves the ability to remember and use the cognitive states on the base level. This is the notion of “conscious” that is opposed to “unconscious” or “subconscious” – and it appears feasible for a multi-layered AI system.

Awareness is often, but not necessarily, connected to a specific way that the cognitive state at the base level *feels* to the subject – this is what philosophers call *phenomenal consciousness*, or how things *seem* to me (Greek *phainetai*). This notion

⁴³ Harry Frankfurt, “Freedom of the will and the concept of a person” *The Journal of Philosophy*, LXVIII: 5; Daniel C. Dennett, *Elbow Room: The Varieties of Free Will Worth Wanting* (Cambridge, Mass. ed, MIT Press, 1984).

⁴⁴ Chapter 6 of this book by Lode Lauwaert and Ann-Katrien Oimann delves further into the subject of AI and responsibility.

⁴⁵ Galen Strawson, “Free will” (2011) *Routledge Encyclopedia of Philosophy*, Taylor and Francis, doi:10.4324/9780415249126-V014-2, www.rep.routledge.com/articles/thematic/free-will/; Thomas Pink, *Free Will: A Very Short Introduction* (Oxford University Press, 2004); Alfred R. Mele, *Free Will and Luck* (Oxford University Press, 2006); Daniel C. Dennett and Gregg D. Caruso, “Just deserts” *Aeon*.

⁴⁶ Thomas W. Simpson and Vincent C. Müller, “Just war and robots’ killings” *The Philosophical Quarterly*, 66: 302; Rob Sparrow, “Killer robots” *Journal of Applied Philosophy*, 24: 62; Vincent C. Müller, “Is it time for robot rights? Moral status in artificial entities” *Ethics & Information Technology*, 23: 579.

of consciousness is probably best explained with the help of two classical philosophical thought experiments: the bat, and the color scientist.

If you and I go out to have the same ice cream, then I can still not know what the ice cream tastes like to you, and I would not know that even if I knew everything about the ice cream, you, your brain, and your taste buds. Somehow, *what it is like* for you is something epistemically inaccessible to me, I can never know it, even if I know everything about the physical world. In the same way, I can never know what it is like to be a bat.⁴⁷

A similar point about what we cannot know in principle is made by Frank Jackson in the article “What Mary didn’t know.”⁴⁸ In his thought experiment, Mary is supposed to be a person who has never seen anything with color in her life, and yet she is a perfect color scientist, she knows everything there is to know about color. One day, she gets out of her black and white environment and sees color for the first time. It appears that she learns something new at that point.

The argument that is suggested here seems to favor an argument for a mental-physical *dualism of substances* or at least *properties*: I can know all the physics, and I cannot know all the phenomenal experience, therefore, phenomenal experience is not part of physics. If dualism is true, then it may appear that we cannot hope to generate phenomenal consciousness with the right physical technology, such as AI. In the form of *substance dualism*, as Descartes and much of religious thought had assumed, dualism is now unpopular since most philosophers assume physicalism, that “everything is physical.”

Various arguments against the reduction of mental to physical *properties* have been brought out, so it is probably fair to say that *property dualism* has a substantial following. This is often combined with substance monism in some version of “supervenience of the mental on the physical,” that is, the thesis that two entities with the same physical properties must have the same mental properties. Some philosophers have challenged this relation between property dualism and the possibility of artificial consciousness. David Chalmers has argued that “the physical structure of the world – the exact distribution of particles, fields, and forces in spacetime – is logically consistent with the absence of consciousness, so the presence of consciousness is a further fact about our world.” Despite this remark, he supports computationalism: “... strong artificial intelligence is true: there is a class of programs such that any implementation of a program in that class is conscious.”⁴⁹

⁴⁷ Thomas Nagel, “What is it like to be a bat?” *Philosophical Review*, 83: 435; Thomas Nagel, *What Does It All Mean? A Very Short Introduction to Philosophy* (Oxford University Press, 1987), chapter 3.

⁴⁸ Frank Jackson, “What Mary didn’t know” *Journal of Philosophy*, 83: 291.

⁴⁹ David J. Chalmers and John R. Searle, “‘Consciousness and the philosophers’: An exchange” (1997) *The New York Review of Books*, www.nybooks.com/articles/1997/05/15/consciousness-and-the-philosophers-an-exchange/; David J. Chalmers, “Précis of the Conscious Mind” (1999) *Philosophy and Phenomenological Research*, LIX(2): 435–438, 436; Donald Davidson, “Mental events” in L. Foster and J. Swanson (eds), *Experience and Theory* (Amherst, MA: University of Massachusetts Press, 1970).

What matters for the function of consciousness in AI or natural agents is not the discussion about dualisms, but rather why phenomenal consciousness in humans is the way it is, how one could tell whether a system is conscious, and whether there could be a human who is physically just like me, but without consciousness (a “philosophical zombie”).⁵⁰

2.8.2 *The Self*

Personal identity in humans is mainly relevant because it is a condition for allocating responsibility (see Section 2.7): In order to allocate blame or praise, there has to be a sense in which I am *the same person* as the one who performed the action in question. We have a sense that there is a life in the past that is mine, and only mine – how this is possible is known as the “persistence question.” The standard criteria for me being the same person as that little boy in the photograph are my *memory* of being that boy, and the *continuity of my body* over time. Humans tend to think that *memory* or *conscious experience*, or *mental content* are the criteria for personal identity, which is why we think we can imagine surviving our death, or living in a different body.⁵¹

So, what is a “part” of that persistent self? Philosophical phantasies and neurological rarities⁵² aside, there is now no doubt what is “part of me” and what is not – I continuously work on maintaining that personal identity by checking that the various senses are in agreement, for example, I try to reach for the door handle, I see my hand touching the handle, I can feel it ... and then I can see the door opening and feel my hand going forward. This is very different from a computer: The components of the standard Von Neumann architecture (input-system, storage, random-access memory, processor, output-system) can be in the same box or miles apart, they can even be split into more components (e.g., some off-board processing of intensive tasks) or stored in spaces such as the “cloud” that are not defined through physical location. And that is only the hardware, the software faces similar issues, so a persistent and delineated self is not an easy task for an AI system. It is not clear that there is a function for a self in AI, which would have repercussions for attributing moral agency and even patiency.

2.9 NORMATIVITY

Let us return briefly to the issues of rational choice and responsibility. Stuart Russell said that “AI has adopted the standard model: we build optimising machines, we

⁵⁰ O'Regan.

⁵¹ Thomas Metzinger, *The Ego Tunnel: The Science of the Mind and the Myth of the Self* (Basic Books, 2009); Eric Olsen, “Personal identity” (2023) *Stanford Encyclopedia of Philosophy*, <https://plato.stanford.edu/entries/identity-personal/>.

⁵² See, for example, “The man who fell out of bed” in Oliver Sacks, *The Man Who Mistook His Wife for a Hat, and Other Clinical Tales* (New York: Summit Books, 1985) or the view of humans as superorganisms, based on the human microbiome.

feed objectives into them, and off they go.”⁵³ On that understanding, AI is a tool, and we need to provide the objectives or goals for it. Artificial intelligence has only *instrumental intelligence* on how to reach given goals. However, *general intelligence* also involves a metacognitive reflection on which goals are relevant to my action now (food or shelter?) and a reflection on which goals one should pursue.⁵⁴ One of the open questions is whether a nonliving system can have “real goals” in the sense required for choice and responsibility, for example, of goals that have subjective value to the system, and that the system recognizes as important after reflection. Without such reflection on goals, AI systems would not be moral agents and there could be no “machine ethics” that deserves the name. Similar considerations apply to other forms of normative reflection, for example, in aesthetics and politics. This discussion in AI philosophy seems to show that there is a function for normative reflection in humans or AI as an elementary part of the cognitive system.

⁵³ Russell, *Human Compatible: Artificial Intelligence and the Problem of Control*, 172.

⁵⁴ Vincent C. Müller and Michael Cannon, “Existential risk from AI and orthogonality: Can we have it both ways?” *Ratio*, 35: 25.