

Can Generative AI Produce Novel Evidence?

Donal Khosrowi (Leibniz University Hannover)

donal.khosrowi@philos.uni-hannover.de

Finola Finn (University of Luxembourg, Leibniz University Hannover)

finola.finn@uni.lu

Abstract

Researchers in archaeology explore the use of generative AI (GenAI) systems for reconstructing destroyed artifacts. This paper poses a novel question: can such GenAI systems generate *evidence* that provides new knowledge about the world or can they only produce *hypotheses* that we might seek evidence *for*? Exploring responses to this question, the paper argues that 1) GenAI outputs can at least be understood as *higher-order evidence* (Parker 2022) and 2) may also produce *de novo synthetic evidence*.

Keywords: Artificial intelligence, Generative AI; Evidence; Historical Science; Archaeology; Computer Simulation; Epistemology

1. Introduction

Artificial intelligence (AI) systems, including generative AI (GenAI), play ever larger roles across the sciences: they are used to make novel discoveries, e.g., of proteins, drugs, or materials (Jumper et al. 2021; Sourati and Evans 2023); to identify new concepts and equations in physics (Iten et al. 2020; Udrescu et al. 2020; Wu and Tegmark 2019); and to suggest new hypotheses, ideas, research questions, or experiments (Krenn et al. 2023; Melnikov et al. 2018). These increasingly extensive roles played by AI put foundational concepts we use to understand and structure scientific pursuits under pressure (Löhr 2023; Hopster and Löhr 2023). For instance, what does it mean to be a scientific ‘discoverer’ (Clark and Khosrowi 2022)? Are AI systems like AlphaFold only ‘tools’ that humans use, or can they exhibit attributes such as ‘autonomy’ or scientific ‘understanding’ (Barman et al. 2024) which we consider essential to the role of ‘discoverer’ or ‘researcher’?

While use cases of AI in physics, chemistry, and biology attract increasing attention by philosophers, there are also underexplored emerging uses of AI in the historical sciences (i.e. history, archaeology, art history, cultural heritage studies, historical anthropology), where researchers explore the use of GenAI for reconstructing partially destroyed manuscripts and artifacts (Navarro et al. 2022; 2023; Lamb et al. 2022; Moral-Andrés et al. 2023). Turning attention to these uses, this paper draws out a novel conceptual disruption regarding how we should understand the *outputs* of GenAI systems: can GenAI systems *generate* evidence that provides genuinely new knowledge in the way that, say, finding new material evidence can? Or can they only produce hypotheses, which may give us reasons for pursuit, but ultimately are the kind of thing that we require evidence *for*? Call this the *evidence question*. Like other conceptual disruptions caused by AI, the evidence question does not have a straightforward answer and highlights substantial uncertainty around how we should apply the concept of ‘evidence’ (see also Rowbottom et al. 2023). The issues this raises are not merely terminological ones but have epistemic and methodological import for practicing researchers. Classifying an output as ‘evidence’ rather than a ‘hypothesis’ confers information about it; in turn, existing norms attached to these classifications may trigger different expectations, attitudes, and actions as appropriate in relation to an output.

Beyond putting the evidence question on the map, this paper also explores potential responses to it. We first consider related debates in the philosophy of computer simulation, where scholars such as Wendy Parker (2022) have elucidated whether simulation systems, such as those used in climate science, can provide (new) evidence for claims about the earth’s climate. Drawing on this debate, we argue that GenAI systems can at least provide *higher-order evidence* in Parker’s sense, i.e. evidence that other evidence for a claim about the world exists. We also explore a more ambitious argument, according to which GenAI systems can produce *de novo synthetic evidence*, which could be epistemically on par with traditional forms of evidence, such as material evidence or expert judgment. The argument suggests they do so by performing pattern recognition-type inferences to yield outputs that provide genuinely new knowledge to agents who lack the ability to make those same inferences. Importantly, while this argument hints at interesting possibilities for understanding GenAI outputs as *de novo synthetic evidence*, it remains agnostic on what historical scientists should or would do with such evidence. In particular, we do not suggest that synthetic evidence is ever an end point or silver bullet for historical and archaeological inquiry (Nygren and Drimmer 2023). If used, it would require description, analysis, contextualisation, and interpretation by historical scientists, as with any other form of evidence.

The discussion is organized as follows. Section 2 outlines the emerging use of GenAI in the historical sciences. Section 3 sharpens the evidence question. Section 4 explores debates in the philosophy of computer simulation and sketches the sequential arguments that GenAI can at least produce higher-order evidence, as well as, possibly, synthetic evidence. Section 5 concretizes why we may take GenAI outputs seriously. Section 6 concludes.

2. Generative AI in the Historical Sciences

A central challenge for researchers in the historical sciences is that the ‘record’ of historical evidence, e.g. manuscripts or artifacts such as pottery, is an imperfect and partial reflection of past events, and is eroded – both figuratively and literally. Not everything survives, and what does is often incomplete or broken. Standard activities in getting a handle on the past (e.g. analyzing the stratigraphic relationships between features at archaeological sites, entertaining larger inferences about chronology, or inferring trade patterns) hence revolve around reconstructing what *was* from what *remains*. Reconstructing partially destroyed artifacts, e.g. to better determine relevant morphological or textural features, is currently often performed by hand, which is resource intensive, can further deteriorate remaining fragments, and cannot deal with fragments that are missing (Navarro et al. 2023). Dealing with these and similar challenges, there is a rich tradition in the historical sciences, especially in archaeology, to recruit technologies from other fields (Wylie 2000), e.g. for sensing and scanning, or, in computational archaeology, using machine-learning methods. For instance, Navarro et al. (2023) develop a GenAI system based on generative adversarial networks (GANs; Goodfellow et al. 2014) called *IberianVoxel*, which reconstructs broken Iberian pottery artifacts as 3D-models. GANs consist of a coupled generator and discriminator architecture; in Navarro et al.’s case, the generator produces 3D-voxel geometries of pottery and the discriminator ‘judges’ whether the geometries produced by the generator look like they were drawn from the data distribution of scanned real artifacts on which it is trained. After a period of adversarial training, the GAN is evaluated, including by surveying domain experts to assess reconstruction quality. The authors report that “archaeologists judge that IberianVoxel generated a correct Iberian style from an initial fragment, and also consider that the reconstructed pottery is between Good and Very Good” (2023, 5839), and conclude their system is “very helpful for exploring and designing automatic procedures to aid experts with the pottery completion task” (ibid., 5833).

Systems such as IberianVoxel are first steps on a trajectory towards more advanced systems permitting finer-grained inferences, especially as GenAI technologies become cheaper to train.

Extrapolating along this trajectory, let us imagine a stylized toy case inspired by IberianVoxel to draw out the central question of this paper more clearly. Consider *AlphaPot*, an imagined GenAI system that has been trained on a very large dataset consisting of images and corresponding high-quality 3D-scans of a wide range of pottery artifacts in various states of decay. AlphaPot is trained to reconstruct masked/corrupted features of an input: i.e. parts of a 3D-model of a scanned real artifact are intentionally corrupted (e.g. by generating synthetic data based on real artifacts that simulate fragmentation, see Lamb et al. 2022, or by physically breaking artifacts and then re-scanning them) and the system is forced to predict how the uncorrupted artifact would have looked like. Assume we are impressed with AlphaPot's performance on unseen test data: it accurately reconstructs broken artifacts for which the ground truth geometry is known. Imagine now that we use AlphaPot to provide a reconstruction R of a novel, partially destroyed artifact A, for which the ground truth is unknown. A is missing pieces that haven't been recovered, but are believed to be essential to classifying A's likely origin or function. R, let us assume, is a plausible looking 3D-model exhibiting fine-grained morphological features that would significantly aid a domain expert (or, for that matter, another AI system) in telling when and where A originated.

3. What's Going on Here, Epistemically?

The key disruption motivating this paper is now clearly in view: what's going on here, epistemically? Has AlphaPot generated a *hypothesis* or made a *prediction*? Or has it generated *evidence*, providing experts with genuinely new knowledge about how A looked like when it was still intact? Understood as a *mere* hypothesis, R is the kind of thing that might give us reasons for pursuit, seeking further evidence to support that R is indeed what A looked like when it was still intact. By contrast, understood as *evidence*, R might already, by itself, support a range of hypotheses regarding A, as well as figure in further downstream inferences that A bears on, e.g. about trade taking place between communities.¹

In making progress on the evidence question, we need to find a benchmark first. A standard Bayesian conception of evidence requires only that a token of evidence E has the capacity to

¹ To be clear, we do not draw a principled distinction between 'hypothesis' and 'evidence.' In line with Bayesian accounts (e.g. Bovens and Hartmann 2003), the difference is contextual. Another way of putting the evidence question is whether R constitutes, or gives rise to, a mere hypothesis H that enjoys no support thus far, and hence has an uninformative or low prior, or whether R constitutes, or gives rise to, a pre-justified hypothesis H' about A that 1) has a high prior and 2) may stand in relevant evidential relationships with yet other hypotheses H'', e.g. about A's likely origin, or whether trade took place between where A was found and other communities.

increase the posterior probability we assign to a hypothesis H (say, a claim about how A looked like) relative to some background theory T (Bovens and Hartmann 2003). It is easy to imagine that R has *some* such capacity, but that is not a very interesting insight (see also Rowbottom et al. 2023 who explore additional complications). Other, functional conceptions of evidence focus on what *role* evidence plays. Here, we are sympathetic to accounts that consider evidence as always being 1) *of* something, 2) *for* something, and 3) *to* someone, relative to a theory of evidence (Hacking 2006; Martini 2021; see also Kosso 2009; Jordanova 2012). For instance, a freshly excavated artifact A is evidence *of* something, e.g. the fact that pottery of A’s kind was made, used, or traded at site S; evidence *for* something, e.g. an inferred claim that pottery of A’s kind was produced in P but ended up at S through a trade route; and evidence *to* someone who has a theory of evidence T and relevant background knowledge K to tell what A can be evidence *of* and *for*. Beyond following such structured conceptions, the subsequent discussion will remain largely uncommitted to specific *philosophical* accounts of evidence. Instead, we find it more productive to consider evidential practices in the historical sciences and think about what existing benchmark types of evidence we could compare GenAI outputs to. What could such benchmarks be? The historical sciences rely on primary sources, e.g. artifacts and documents that are close (causally, spatially, temporally, by provenance) to the phenomena of interest. Relevant benchmarks to address the evidence question could hence be, for instance, a highly similar, intact artifact B found in the same stratum at the same site, or pertinent text, illustrations, or tools bearing on the likely morphological features of A. Likewise, expert judgment that joins up available background theory and primary material evidence in a larger inference is another candidate. The evidence question is sharper now: could AlphaPot’s outputs be considered evidence comparable to these benchmarks, e.g. other, material evidence like B that could licence an analogical inference that ‘A would have probably looked like B when it was intact’ or expert judgment that joins various resources together to yield, say, a rendition or description of what A would look like, were it still intact?

4. Yes, but What Kind of Evidence Is It?

The answer that we want to explore here is: yes, GenAI systems like AlphaPot have the capacity to generate *synthetic evidence* that provides genuinely new knowledge about the world. What could an argument for such a thesis look like? A first pass could build on familiar successes of using AI systems for inferential tasks in science, like AlphaFold 2.0 (Jumper et al. 2021). Specifically, at training, these systems 1) latch onto information, especially high-

dimensional and distributed correlational information or patterns, in training data, and 2) learn a model, i.e. an abstract representational space encoding relevant features and a corresponding function $f(\cdot)$ within that space, which maps inputs to outputs in a way that minimizes empirical risk (at least locally). At inference, such models, given an input (e.g. a scan of a partially destroyed artifact), 3) generate outputs that yield accurate reconstructions of the input, as governed by $f(\cdot)$. This kind of story could touch on guarantees in machine learning (e.g. Cybenko 1989) and statistical learning theory (e.g. Vapnik 2000; Bargagli Stoffi et al. 2022) to explain notable successes of machine learning systems, e.g. in latching onto complex, subtle, and distributed patterns that escape human attention, such as in skin cancer classification or protein structure prediction (Jumper et al. 2021), or in successfully learning novel high-dimensional representations (e.g. word or image embeddings) that can be used for text and image synthesis, as demonstrated by GenAI systems like ChatGPT or StableDiffusion (Rombach et al. 2022).

This story, while somewhat compelling, is still too simple. Here, we focus on concerns arising from related debates in the philosophy of scientific models and computer simulation. In this space, philosophers have tried to understand whether models and simulations can provide genuinely new knowledge about the world and, if so, how (e.g. Parker 2022; Beisbart 2012). In a nutshell, sceptics about the epistemological significance of models and simulations point out that these tools only help us recognize the consequences of knowledge that we already possess, e.g. assumptions (e.g. equations) and initial conditions (e.g. measurements, parameterizations). These consequences can at most be evidence in the sense that they provide new information to agents who are not able to, or simply did not, derive those same consequences given the same assumptions and initial conditions. But they would not be evidence to a more ideal agent who would already recognize these consequences under some form of inferential closure. So, while observation and experimentation allow us to gather *new* experience (Beisbart 2012, 245), models and simulations don't bring anything new to the table; though they do help limited epistemic agents better see what's already on the table.

What does this mean for the evidence question? Parker summarizes the consequences of the sceptical view on computer simulation as follows: “If computer simulation is at bottom an attempt to calculate the implications of a set of modelling assumptions, then simulation results ... seem to be *predictions rather than evidence*; they are the kind of thing we might seek evidence for”. (Parker 2022, 1522; emphasis added). On such a view, the outputs of GenAI systems like AlphaPot are predictions, or, more generally, hypotheses. They might alert us to possibilities for how an artifact may have looked like, and may give us reasons for

pursuing these hypotheses by means of bringing evidence to bear on them; but they are not to be taken as evidence that could already, by itself, support knowledge claims about artefacts or figure importantly alongside other evidence in larger, downstream inferences, such as about trade taking place between different communities.

Filling the space between more extreme views that either consider simulation results evidence, or deny that they can ever be, Parker offers a finer-grained view to characterize what simulations provide to agents. Specifically, Parker argues that simulation outputs can be *higher-order evidence*: they can be evidence E that other evidence E' for a hypothesis H exists. Specifically, such higher-order evidence can help agents obtain genuinely new knowledge of the world if they 1) either don't have access to E', or else, 2) lack the background knowledge needed to understand how E' bears on H. So, while simulations “do not provide information about the world that goes beyond that which is already implicit in their assumptions, particular epistemic agents—including even scientists and engineers using simulation models—might still gain genuinely new knowledge of the world via simulation.” (Parker 2022, 1522)

Parker's view offers a useful backstop for thinking about GenAI outputs. At the very least, they seem able to figure as higher-order evidence. A 3D reconstruction of a broken artifact A from a suitably validated system like AlphaPot provides new knowledge about specific artifacts to agents who either don't have access to the training data² E' that bear on the reconstructive query about A, or else lack the background knowledge to understand how E' bears on questions about A. This is a useful insight already, but it also seems interesting to explore whether GenAI systems could ever provide more than 'just' higher-order evidence.

4.1 More Than Higher-Order?

What might GenAI systems be doing that goes beyond what simulation systems do? A central difference seems to be that GenAI systems can exhibit higher degrees of independence, which allows them to perform computations that instantiate *inferences* of a kind that simulation systems don't instantiate. Specifically, simulation systems in the climate sciences are built based on highly developed antecedent understanding of the physics equations describing aspects of the earth's climate system (background knowledge), parameterized according to our best understanding of key parameters and known/understood aspects of the phenomena

² In particular, they might lack access to the *information* contained in that data, e.g. regularities about the 'grammar' of Iberian pottery.

involved, and calibrated using data regarding the earth’s climate system. Together, these inputs substantially constrain the behaviors of simulating systems.

GenAI systems exhibit comparatively higher independence because they are not as tightly constrained. There are no accepted equations that describe, say, the ‘grammar’ of Iberian pottery. Nor, for lack of such equations, are there measurements that GenAI systems are parameterized with. In short, there is no developed body of background knowledge that is explicitly encoded when building GenAI systems (at least in unsupervised/self-supervised regimes), nor would our existing background knowledge permit building systems in a way that mirrors the strategies behind building simulation systems. Rather, the very purpose of machine learning approaches is often to *extract* pertinent background knowledge from data, e.g. to find a function $f(\cdot)$ that usefully captures features of a joint distribution and can be used to perform successful inferences. For this enterprise to be successful, GenAI systems must exhibit considerable degrees of freedom to ‘settle’ on representational spaces, representations, and input-output relationships that are 1) predictively useful, 2) possibly inaccessible to humans by other means (e.g. visual inspection), and 3) potentially novel to humans. GenAI systems hence harbor the capacity for a special kind of novelty in their outputs. Unlike simulation systems, they can generate *synthetic evidence*, i.e. evidence E that is not only psychologically novel to agents who lack other evidence E’ or background knowledge K, but is novel to agents who do not possess the same *inferential abilities* to extract pertinent knowledge K (e.g. of $f(\cdot)$) from the same training data. Such abilities are different from *computational* abilities to derive *implications* of equations and initial conditions. They are more akin to the ability to ‘recognize’ that such-and-such is a good way to represent or compress data, or that such-and-such is a successful (i.e. error-minimizing) way to ‘fill in the blanks’ of a reconstructive query.

On the narrative presented here, GenAI systems bring inferential abilities to the table that simulation systems don’t. But why should this lead us to conclude that they can *generate* evidence that provides genuinely new knowledge to agents? Couldn’t, or shouldn’t, we still maintain that the relevant information with bearing on reconstructive queries ‘resides in’ the training data that GenAI systems are trained on?³ This would bring us back to understanding GenAI systems as, at most, providing higher-order evidence in Parker’s sense and conclude that no evidence that is novel over and above whatever is contained in these data is generated.

A good way to explore how GenAI systems can provide novelty beyond higher-order evidence is to think about patterns. A standard success narrative of machine-learning based inference

³ This concern also flags a version of the *problem of old evidence*; see e.g. Sprenger (2015) for a discussion.

alluded to above is centrally tied to systems’ abilities to identify patterns, including subtle and distributed ones, and to exploit them for inference. But what is a pattern, anyway? This is yet another issue that ML systems press us to confront with greater care. Here, it is useful to distinguish two general types of views on patterns: *ontic* and *epistemic* views. On the first, a pattern is constituted by a collection of material facts about the world that may be distributed across entities. A pattern, on this view, is always ‘there’, even without a mind to recognize and exploit it for inference (see Ladyman and Ross 2007 on ‘real patterns’). On an epistemic view (cf. Dennett 1991; McAllister 2010; Haugeland 1998; Kästner and Haueis 2021), patterns come into existence through an epistemic *agent* that recovers it, including, say, by devising an ontology of entities and features ranging over a domain (e.g. pots, fractures, materials, textures); making efforts to describe and represent these entities and features within that domain in an abstract way (e.g. material or shape types); and exploring how these representations hang together, e.g. causally or probabilistically, at that abstract representational level. On such a view, a pattern is instantiated by, refers to, and supervenes on, concrete material things, but ultimately resides at an abstract representational level (cf. Ladyman and Ross 2007 on ‘second-order patterns’). If we find such an epistemic view compelling, then this allows that GenAI systems, like other epistemic agents, can perform inferential activities that *bring patterns into existence*.⁴ This ability sets GenAI systems apart from simulating systems: they may produce outputs, based on the ability to infer patterns from data, that are novel to agents who do not possess such abilities.⁵

5. Strictures on Synthetic Evidence

We now have a sketch of an argument for the claim that GenAI systems like AlphaPot may produce *synthetic evidence*, i.e. evidence E that provides genuinely new knowledge about the world to agents who do not possess the same inferential abilities to recover E from primary evidence E’ as the system that produced E. But when can we expect GenAI systems to produce *good* synthetic evidence? As researchers are exploring use-cases of LLMs in history, for instance to ‘ventriloquize’ the voices of the past through LLMs trained and/or fine-tuned on historical text corpora to enable researchers to ‘query’ past societies or individuals (Hutson et al. 2024), there is a real risk of low-cost bogus AI-driven science. While there are a variety

⁴ This view is still compatible with realist views like Ladyman and Ross’ (2007); Epistemic agents, or GenAI systems on our account, bring into existence real second-order patterns that represent real first-order patterns.

⁵ Of course, we must mind anthropomorphic pitfalls. Terms like ‘recognizing’, ‘using’, and so on, must not be taken to suggest that GenAI systems literally have mental states or cognitive abilities associated with these terms.

of salient concerns about the reliability of GenAI systems, such as regarding ‘hallucinations’, brittleness, lack of generalization abilities, and epistemic opacity, here, we outline some potential virtues that GenAI systems may exhibit, if designed and deployed responsibly. These virtues help better understand the conditions under which we may reasonably hope these systems to make valuable epistemic contributions.

- 1) **Scope:** GenAI systems are good at processing and ‘drawing on’ large amounts of rich data, which is relevant when patterns are distributed across large numbers of entities and different data modalities.
- 2) **Sensitivity:** ML systems are known to usefully latch onto subtle, distributed patterns, especially in quantitative data, that are often not accessible to human perception.
- 3) **Probabilism:** ML-based inference is probabilistic. Outputs are sampled from a whole modeled joint distribution. This often means that other possibilities for an output are not discarded by a system, but remain, or could be made, available to investigators.
- 4) **Mechanicity:** outputs are often (near-) repeatable from the same inputs, so GenAI systems can be subjected to systematic intervention, allowing investigators to understand how outputs depend on inputs. For instance, they may upsample rare input types (e.g., by using synthetic training data to induce more variation regarding specific artifact types) and gauge whether outputs change for specific query types.
- 5) **Theory-freedom/-agnosticism:** especially in unsupervised or self-supervised learning regimes, GenAI systems organize data somewhat independently of existing theory, working against unhelpful forms of theory-laden observation.
- 6) **Complexity:** universal function approximation theorems (e.g. Cybenko 1989) and statistical learning theory (Vapnik 2000; Bargagli Stoffi et al. 2022) provide (probabilistic) guarantees for specific system-types to successfully approximate arbitrarily complex input-output relationships under suitable conditions. This is important as there are no good reasons to believe that, say, the ‘grammar’ of Iberian pottery (i.e. the ‘rules’ that govern the joint distribution of morphological features of Iberian pottery artifacts) is easily captured by simple, human-expressible functions.
- 7) **Granularity:** GenAI systems perform inference at multiple levels, including at fine-grained pixel- or voxel-levels that may not be salient to human investigators. Such systems are hence not as susceptible as humans to latch exclusively onto patterns or analogies obtaining at higher, more salient levels of analysis, e.g. inferring that artifact A probably had inscription S because, B, C, D, who look morphologically similar, do.

Of course, specific GenAI systems are not guaranteed to exhibit any of these virtues to significant degrees; only well-engineered systems may. Moreover, many of the candidate virtues outlined here can turn into vices if the properties they track are expressed too strongly: think of theory-freedom or scope that could lead a system to consider irrelevant or misleading information when there are good theoretical reasons not to. Spelling out a contextualist virtue epistemology for GenAI systems in science is arguably a larger project that will require more space, which is why the virtues sketched here should only provide some early inspiration rather than a sketch of a full-fledged account of what GenAI systems may bring to the table. That said, it seems promising to explore such an account in articulating answers to the evidence question.

6. Conclusions

This paper puts an important new question about the role of generative AI (GenAI) systems in the sciences on the map. The *evidence question* asks: can GenAI systems *generate* evidence that provides agents, including experts, with genuinely new knowledge about the world? Focusing on the historical sciences, where researchers explore the use of GenAI systems to reconstruct partially destroyed manuscripts and artifacts to learn about the past, we argued that it is currently unclear whether we should understand the outputs produced by these systems as mere *hypotheses* or as *evidence*, where the former may give researchers reasons for pursuit and for seeking out further evidence, and the latter may already licence knowledge claims about the world and figure directly in supporting further inferences. Given this conceptual and practical uncertainty, we sketched how we may understand GenAI outputs not only as *higher-order evidence* in the sense of Parker (2022) but also potentially as *synthetic evidence*, i.e. evidence that can provide agents, including experts, with genuinely new knowledge about the world. They do so by acquiring and deploying pattern recognition-type inferential abilities to produce outputs that are evidence to agents who lack those same inferential abilities, which may include even our best domain experts. The scope of this argument sketch is narrow: it applies, for now, only to the emerging uses of GenAI in the historical sciences discussed here. But zooming out, the evidence question may also extend to a range of other domains that explore the utility of GenAI for advanced inferential tasks (e.g. drug and materials discovery). For philosophers of science this is good news, inviting us to help characterize and resolve the methodological disruptions affecting emerging scientific practices, and to contribute to development of sound methodologies involving GenAI.

References

- Barman, Kristian G., Sascha Caron, Tom Claassen, and Henk de Regt. 2024. “Towards a Benchmark for Scientific Understanding in Humans and Machines”. *Minds & Machines* 34:6. <https://doi.org/10.1007/s11023-024-09657-1>
- Bargagli Stoffi, Falco J., Gustavo Cevolani, and Giorgio Gnecco. 2022. “Simple Models in Complex Worlds: Occam’s Razor and Statistical Learning Theory.” *Minds & Machines* 32:13–42. <https://doi.org/10.1007/s11023-022-09592-z>
- Beisbart, Claus. 2012. “How Can Computer Simulations Produce New Knowledge?” *European Journal for Philosophy of Science* 2:395–434. <https://doi.org/10.1007/s13194-012-0049-7>
- Bovens, Luc, and Stephan Hartmann. 2003. *Bayesian Epistemology*. Oxford: Oxford University Press. <https://doi.org/10.1093/0199269750.001.0001>
- Clark, Elinor, and Donal Khosrowi. 2022. “Decentering the Discoverer: How AI Helps Us Rethink Scientific Discovery.” *Synthese* 200:463. <https://doi.org/10.1007/s11229-022-03902-9>
- Cybenko, George. 1989. “Approximation by Superpositions of a Sigmoidal Function.” *Mathematics of Control, Signals, and Systems* 2:303–14. <https://doi.org/10.1007/BF02551274>
- Dennett, Daniel C. 1991. “Real Patterns.” *Journal of Philosophy* LXXXVIII:27–51.
- Durán, Juan M., and Formanek, Nico. 2018. “Grounds for Trust: Essential Epistemic Opacity and Computational Reliabilism.” *Minds & Machines* 28:645–66. <https://doi.org/10.1007/s11023-018-9481-6>
- Goodfellow, Ian, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. “Generative Adversarial Nets.” *Advances in Neural Information Processing Systems* 27.
- Hacking, Ian. 2006. *The Emergence of Probability: A Philosophical Study of Early Ideas About Probability, Induction and Statistical Inference*. Cambridge: Cambridge University Press. <https://doi.org/10.1017/CBO9780511817557>
- Haugeland, John. 1998. “Pattern and Being.” In *Having Thought: Essays in the Metaphysics of Mind*, 267–90. Cambridge, MA: Harvard University Press.
- Hopster, Jeroen, Philip Brey, Michael Klenk, Guido Löhr, Samuela Marchiori, Björn Lundgren, and Kevin Scharp. 2023. “Conceptual Disruption and the Ethics of Technology.” In *Ethics of Socially Disruptive Technologies: An Introduction*, ed. Ibo van de Poel, Lily Frank, Julia Hermann, Jeroen Hopster, Dominic Lenzi, Sven Nyholm,

- Behnam Taebi, and Elena Ziliotti, 141-62. Cambridge, UK: Open Book Publishers.
<https://doi.org/10.11647/OBP.0366.06>
- Hopster, Jeroen, and Guido Löhr. 2023. “Conceptual Engineering and Philosophy of Technology: Amelioration or Adaptation?” *Philosophy & Technology* 36:70.
<https://doi.org/10.1007/s13347-023-00670-3>
- Hutson, James, Paul Huffman, and Jeremiah Ratican. 2024. “Digital Resurrection of Historical Figures: A Case Study on Mary Sibley through Customized ChatGPT.” *Metaverse* 4 (2): 2424. <https://doi.org/10.54517/m.v4i2.2424>
- Iten, Raban, Tony Metger, Henrik Wilming, Lidia del Rio, and Renato Renner. 2020. “Discovering Physical Concepts with Neural Networks.” *Physical Review Letters* 124:010508. <https://doi.org/10.1103/PhysRevLett.124.010508>
- Jiang, Harry H., Lauren Brown, Jessica Cheng, Mehtab Khan, Abhishek Gupta, Deja Workman, Alex Hanna, Johnathan Flowers, and Timnit Gebru. 2023. “AI Art and its Impact on Artists.” In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society (AIES '23)*, 363–74. <https://doi.org/10.1145/3600211.3604681>
- Jordanova, Ludmilla. 2012. *The Look of the Past: Visual and Material Evidence in Historical Practice*. Cambridge: Cambridge University Press.
<https://doi.org/10.1017/9781139051095>
- Jumper, John, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstern, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. 2021. “Highly Accurate Protein Structure Prediction with AlphaFold.” *Nature* 596:583-89. <https://doi.org/10.1038/s41586-021-03819-2>
- Kästner, Lena, and Philipp Haueis. 2021. “Discovering Patterns: On the Norms of Mechanistic Inquiry.” *Erkenntnis* 86:1635-60. <https://doi.org/10.1007/s10670-019-00174-7>
- Kosso, Peter. 2009. “Philosophy of Historiography.” In *A Companion to the Philosophy of History and Historiography*, ed. Aviezer Tucker, 7-25. Chichester: Wiley Blackwell.
<https://doi.org/10.1002/9781444304916>
- Krenn, Mario, Lorenzo Buffoni, Bruno Coutinho, Sagi Eppel, Jacob Gates Foster, Andrew Gritsevskiy, Harlin Lee, Yichao Lu, João P. Moutinho, Nima Sanjabi, Rishi Sonthalia, Ngoc Mai Tran, Francisco Valente, Yangxinyu Xie, Rose Yu, and Michael Kopp. 2023. “Forecasting the Future of Artificial Intelligence with Machine Learning-Based Link

- Prediction in an Exponentially Growing Knowledge Network.” *Nature Machine Intelligence* 5:1326-35. <https://doi.org/10.1038/s42256-023-00735-0>
- Ladyman, James, Don Ross, and David Spurrett with John Collier. 2007. *Every Thing Must Go: Metaphysics Naturalized*. Oxford: Oxford University Press. <https://doi.org/10.1093/acprof:oso/9780199276196.001.0001>
- Lamb, Nikolas, Sean Banerjee, and Natasha Kholgade Banerjee. 2022. “DeepJoin: Learning a Joint Occupancy, Signed Distance, and Normal Field Function for Shape Repair.” *ACM Transactions on Graphics (TOG)* 41 (6): 230. <https://doi.org/10.1145/3550454.3555470>
- Löhr, Guido. 2023. “Conceptual Disruption and 21st Century Technologies: A Framework.” *Technology in Society* 74:102327. <https://doi.org/10.1016/j.techsoc.2023.102327>
- Martini, Carlo. 2021. “What ‘Evidence’ in Evidence-Based Medicine?” *Topoi* 40:299–305. <https://doi.org/10.1007/s11245-020-09703-4>
- McAllister, James W. 2010. “The Ontology of Patterns in Empirical Data.” *Philosophy of Science* 77 (5): 804–14. <https://doi.org/10.1086/656555>
- Melnikov, Alexey. A., Hendrik P. Nautrup, Mario Krenn, and Hans H. Briegel. 2018. “Active Learning Machine Learns to Create New Quantum Experiments.” *PNAS* 115 (6): 1221-26. <https://doi.org/10.1073/pnas.1714936115>
- Moral-Andrés, Fernando, Elena Merino-Gómez, Pedro Reviriego, and Fabrizio Lombardi. 2023. “Can Artificial Intelligence Reconstruct Ancient Mosaics?” *Studies in Conservation* 69 (5): 313–26. <https://doi.org/10.1080/00393630.2023.2227798>
- Navarro, Pablo, Celia Cintas, Manuel Lucena, José Manuel Fuertes, Rafael Segura, Claudio Delrieux, and Rolando González-José. 2022. “Reconstruction of Iberian Ceramic Potteries Using Generative Adversarial Networks.” *Scientific Reports* 12:10644. <https://doi.org/10.1038/s41598-022-14910-7>
- Navarro, Pablo, Celia Cintas, Manuel Lucena, José Manuel Fuertes, Antonio Rueda, Rafael Segura, Carlos Ogayar-Anguita, Rolando González-José, and Claudio Delrieux. 2023. “IberianVoxel: Automatic Completion of Iberian Ceramics for Cultural Heritage Studies.” In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence (IJCAI '23)*, 5833–41. <https://doi.org/10.24963/ijcai.2023/647>
- Nygren, Christopher, and Sonja Drimmer. 2023. “Art History and AI: Ten Axioms”. *International Journal for Digital Art History* 9:5.02-5.13. <https://doi.org/10.11588/dah.2023.9.90400>
- Parker, Wendy S. 2022. “Evidence and Knowledge from Computer Simulation.” *Erkenntnis* 87:1521–38. <https://doi.org/10.1007/s10670-020-00260-1>

- Rombach, Robin, Andreas Blattmann, Dominik Lorenz, Patrick Esser and Björn Ommer. 2022. “High-Resolution Image Synthesis with Latent Diffusion Models.” arXiv. <https://doi.org/10.48550/arXiv.2112.10752>
- Rowbottom, Darrell P., André Curtis-Trudel, and William Peden. 2023. “Evidence, Computation and AI: Why Evidence Is Not Just in the Head.” *Asian Journal of Philosophy* 2: 11. <https://doi.org/10.1007/s44204-023-00061-7>
- Sourati, Jamshid, and James A. Evans. 2023. “Accelerating Science with Human Aware Artificial Intelligence.” *Nature Human Behaviour* 7:1682-96. <https://doi.org/10.1038/s41562-023-01648-z>
- Sprenger, Jan 2015. “A New Solution to the Problem of Old Evidence.” *Philosophy of Science* 82 (3): 383–401.
- Udrescu, Silviu-Marian, Andrew Tan, Jiahai Feng, Orisvaldo Neto, Tailin Wu, and Max Tegmark. 2020. “AI Feynman 2.0: Pareto-Optimal Symbolic Regression Exploiting Graph Modularity.” *Advances in Neural Information Processing Systems 33 (NeurIPS 2020)*, 4860–71.
- Vapnik, Vladimir N. 2000. *The Nature of Statistical Learning Theory*. New York: Springer.
- Wu, Tailin, and Max Tegmark. 2019. “Toward an Artificial Intelligence Physicist for Unsupervised Learning”. *Physical Review E* 100 (3): 033311. <https://doi.org/10.1103/PhysRevE.100.033311>
- Wylie, Alison. 2000. “Questions of Evidence, Legitimacy, and the (Dis)Unity of Science.” *American Antiquity* 65 (2): 227-37. <https://doi.org/10.2307/2694057>

Acknowledgements

We thank the audiences at PSA 2024, Aarhus University, DKPhil2024, the University of Groningen, the GRK-SOCRATES Colloquium, the Machine Discovery and Creation Virtual Workshop, LICPOS 2023, and the 2023 MCMP-LUH-Wuppertal Workshop for their many helpful questions and suggestions, which contributed significantly to refining the arguments presented here.

Funding Information

The research for this article was supported by a grant from the Ministry of Science and Culture of Lower Saxony (MWK), Grant No.: 11-7620-1155/2021.