# Apriori Knowledge in an Era of Computational Opacity: The Role of AI in Mathematical Discovery

Eamon Duede*[1] and Kevin Davey†[2]

[1]Harvard University
[2]University of Chicago

**Abstract**

Computation is central to contemporary mathematics. Many accept that we can acquire genuine mathematical knowledge of the Four Color Theorem from Appel and Haken's program insofar as it is simply a repetitive application of human forms of mathematical reasoning. Modern LLMs / DNNs are, by contrast, opaque to us in significant ways, and this creates obstacles in obtaining mathematical knowledge from them. We argue, however, that if a proof-checker automating human forms of proof-checking is attached to such machines, then we can obtain apriori mathematical knowledge from them, even though the original machines are entirely opaque to us and the proofs they output are not human-surveyable.

## 1   Introduction

The recent emergence of deep learning and generative large language models as tools in pure mathematics [DVB+21, RPBN+23, TWL+24] creates an opportunity to connect current debates in the philosophy of AI-infused science with earlier debates in the philosophy of computation-infused mathematics.

Earlier debates about computation-infused mathematics centered on Appel and Haken's famous 1977 computer proof of the Four Color Theorem [AH89]. Because the computations performed by Appel and Haken's computer are not directly human-checkable due to their length, it appears that our belief in the Four Color Theorem critically depends on our inductive confidence in the reliability of computers [Tym79]. Some thought that this meant that the concept of mathematical knowledge had to be expanded to allow for 'experimental' knowledge [Tym98], and that the idea of mathematical knowledge as essentially apriori had to be abandoned [DL80].

However, following the work of Burge [Bur98], we argue in Section 2 that so long as the running of a computer program can be understood as a mechanized exercise of something like ordinary human mathematical capacities, the output of a program can indeed give us apriori

*Email: eduede@g.harvard.edu
†Email: kjdavey@uchicago.edu

mathematical knowledge. More specifically, we follow Burge in claiming that Appel and Haken did indeed attain apriori knowledge of the truth of the Four Color Theorem from the output of their computer program.

The problem, however, is that this sort of argument does not apply to the output of machines like deep neural networks (henceforth DNNs) and generative language learning models (henceforth LLMs), whose inner workings are, in a sense, opaque to us. In Section 3, we argue that outside special cases, in general, we cannot *directly* acquire mathematical knowledge from the reports of DNNs or LLMs. A result of this kind seems to impose a strong limitation on our ability to acquire genuine mathematical knowledge from artificial intelligence.

Nevertheless, in Section 4, we argue that mathematicians can overcome this limitation by applying a transparent proof-checker to an appropriately structured output of a DNN or LLM. So long as this proof-checker may be understood as a mechanized exercise of human proof-checking capacities, we claim that we can acquire genuine mathematical knowledge using opaque DNNs or LLMs from the output of the proof-checker, even though this knowledge may not be obtained directly from the DNN or LLM itself.

In this way, we arrive at the perhaps surprising result that it is possible to acquire genuine mathematical knowledge of a fact whose proof was generated by a process that is not mathematically intelligible to us, *and* is so complex that it is not human-checkable in any way. This suggests that AI can indeed play a significant and potentially transformative role in generating genuine mathematical knowledge in pure mathematics.

## 2   Knowledge of the Four Color Theorem

The question of what effect computer proofs have on our philosophical conception of mathematics began to be discussed in earnest in 1977 when Appel and Haken used a computer to verify the Four Color Theorem [AH89]. Appel and Haken reduced the Four Color Theorem to that of verifying that a certain 'reducibility' property held of each of a particular set of 1,834 finite graphs. While verifying that a finite graph is reducible can be done mechanically, larger graphs require the consideration of a finite but extremely large number of possibilities. Appel and Haken used a supercomputer to verify that each of these 1,834 graphs is indeed reducible. This required over a month of continual computer operation. As a result, the successful execution of their algorithm can be understood as a proof of the Four Color Theorem which, in virtue of its length and complexity, has yet to be (and presumably cannot be) surveyed and checked step-by-step by a human mathematician. To this day, no proof of the Four Color Theorem exists that can be surveyed and checked by a human.

As early as 1979 philosophers began to reflect on what this meant for mathematics. Tymoczko [Tym79] argued that mathematics had now become an empirical discipline insofar as proofs could be obtained by performing *experiments* such as the running of computer programs. Detlefsen and Luker [DL80] argued that mathematics had in some sense always been an empirical discipline and that there was therefore nothing philosophically new in Appel and Haken's accomplishment.

The philosophical position that we find more attractive and on which we wish to dwell is

a rather different one due to Burge [Bur98]. Rejecting the views just described, Burge argued that the output of Appel and Haken's program gives us an *apriori* (and not merely empirical) warrant for believing the Four Color Theorem.[1] Mathematical justification thus retains its apriori character despite the essential involvement of computers.

In calling a justification or warrant *apriori*, we mean that it does not depend upon empirical considerations in any way for its force. Although the question of how to precisely characterize the apriori is vexed (see for example [Wil13]), in what follows we only need to rely on the fact that traditional and ordinary forms of mathematical argument are apriori, and will not need to posit anything controversial about the nature of the apriori.

Crucial to Burge's argument that we can obtain apriori knowledge from the report of Appel and Haken's computer is the fact that the kinds of capacities that the computer exercises in verifying the Four Color Theorem are the same sorts of rational capacities used by human mathematicians. There is after all some sense in which the way the computer mechanically verifies of a particular graph that it is reducible is precisely the same as the way that a human mathematician would, albeit more quickly and indefatigably. At each step of the computer's operation, Appel and Haken could thus correctly describe what the computer is doing by saying, for example, 'the computer is now considering graph number 1,002, and is verifying that a 4-coloring of a particular subgraph is extendible to the entire graph, by doing such-and-such.' Appel and Haken recognize what the machine is doing as nothing other than a mechanized mobilization and application of ordinary human mathematical capacities. We shall capture all this by saying that the operation of Appel and Haken's program is *mathematically transparent* to them.[2]

Burge goes on to argue that *'we have apriori prima facie entitlement to accept intelligible presentations-as-true, expressed by the print-outs'* [Bur98, p.13] and that *'resources for rationality are, other things equal, to be believed'* [Bur98, p.5]. The idea here is that we have a default, apriori entitlement to accept as correct the output of mathematically transparent processes and that Appel and Haken, knowing that their computer program is mathematically transparent, thus have an apriori warrant for believing the Four Color Theorem. Suppose one further believes, as Burge does, that apriori warrant can be transferred testimonially to others in the right circumstances. In that case, we too have an apriori warrant for believing the Four Color Theorem when we read Appel and Haken's report claiming that a program mathematically transparent to them has returned the output in question.

Burge emphasizes that while apriori, this warrant for believing the Four Color Theorem is defeasible. For example, Appel and Haken could come to learn that the computer was malfunctioning, in which case they would no longer be justified in believing the Four Color Theorem (without further reasons.) On Burge's view, apriori warrants are not generally infallible and can be defeated. The crucial claim, however, is that this does *not* mean that being justified in believing the Four Color Theorem first requires that we have positive empirical grounds for

---

[1]Of course, the output of Appel and Haken's program only assures us that the given 1,834 graphs are reducible, and to get the Four Color Theorem from this we need to supplement it with a further piece of human-generated mathematics. For the sake of compactness however we shall be slightly sloppy and talk of the program as giving us an apriori warrant for believing the Four Color Theorem, even though technically we (and Burge) should only claim that it gives us an apriori warrant for believing that the given 1,834 graphs are reducible.

[2]This is not Burge's exact terminology, but we take it to be close to the spirit of what he says.

thinking that the computer was not malfunctioning. Instead, we are entitled to believe the results of rationally transparent processes so long as we *lack* reason for thinking the source unreliable. So, it is enough that we have no reason to think that the computer was malfunctioning. Appel and Haken's ultimate ground for accepting the Four Color Theorem is thus simply the existence of a rationally transparent process demonstrating it. This warrant is apriori insofar as it derives its force from nothing other than purely mathematical considerations (even though, in coming to know it, we may invoke the reliability of computers and of our own faculties.)

In fact, nothing here is strikingly different from the way in which we treat human mathematicians. In ordinary circumstances when I devise or read a proof of $X$, I acquire an apriori warrant for believing $X$. The force of the warrant is apriori because it involves purely mathematical, non-empirical considerations. This warrant is apriori even though I must rely on my memory, eyesight, and any other number of cognitive capacities in devising or reading the proof, and in spite of the fact that the warrant in question can be defeated by pointing out that I have misremembered or misread something.

If there is a difference here, it is only that Appel and Haken are not directly performing the ordinary mathematical reasoning but have outsourced it to a computer instead. This means, to use Burge's terminology, that the fine details of the argument for the Four Color Theorem fall outside their 'proprietary warrant' – that is, outside the set of reasons directly available to them – though they fall inside their 'extended warrant' – that is, inside the set of facts about the mathematically transparent process on which they depend for their knowledge. The fact that this extended warrant consists of ordinary mathematical considerations and not empirical considerations means that the warrant remains an apriori one.

For a useful contrast, Burge considers the case in which we encounter a mathematical genius whose explanations are so opaque that we have no real mathematical understanding of their reasoning. Insofar as the reasoning of the genius is not rationally transparent to us, Burge thinks that we do not have the type of apriori warrant for believing them that Appel and Haken have for believing the output of their computer. Of course, if the genius's track record is sufficiently strong, we can have an *inductive* warrant for believing them. We do not think however that this amounts to mathematical knowledge. Because of the absense of mathematical transparency, we therefore do not acquire mathematical knowledge from the reports of the genius.[3]

In sum, when a mathematically transparent machine outputs a mathematical claim $X$, we acquire an apriori warrant for believing $X$, and in appropriate circumstances, even obtain apriori knowledge of $X$. In precisely this manner, we take Appel and Haken to have acquired apriori knowledge of the Four Color Theorem.

## 3  Transparency and AI Assisted Proof

More recently, mathematicians have turned to deep learning models (DLMs) for assistance with particularly challenging mathematical problems in, for instance, low-dimensional topology [DVB$^+$21], geometry [TWL$^+$24], and extremal combinatorics[RPBN$^+$23]. If it is the case that apriori warrants for mathematical knowledge of the kind involving computers discussed above

---

[3]It is not entirely clear whether Burge would agree with this last point.

require mathematically transparent computational procedures, then the notorious opacity of deep learning approaches presents a straightforward problem for the possibility of getting genuine mathematical knowledge from them.

To see this, we need to consider the sense in which DLMs are opaque and, as such, cannot be made mathematically transparent to us. Creel's account of algorithmic and structural transparency in complex computational systems [Cre20] is particularly helpful in this regard and, as we will show, can serve to connect long-running debates concerning computational approaches in mathematics with emerging contemporary debates in the philosophy of AI-infused science.

For Creel, computational systems can be 'algorithmically' and 'structurally' transparent.[4] A computational system is said to be algorithmically transparent to the extent that the procedures that govern its behavior are known and intelligible. In the case of the procedures performed in the proof of the Four Color Theorem, the rules at the algorithmic level correspond to the straightforward application of well-accepted principles of graph theory. The system is structurally transparent to the extent that it is possible to see how the mathematical principles are realized in code. Thus, we can say that the program is structurally transparent if and only if it is both surveyable and possible to understand how the code generates results in accordance with the general mathematical principles that it instantiates.[5]

While the computations used in the Four Color Theorem resulted in an unsurveyable *proof*, the computations themselves were, on Creel's account, fully algorithmically and structurally transparent. We claim this means that the computations are mathematically transparent in the sense discussed in the last section. However, we will argue that the use of both DNNs and LLMs in mathematics often involves dependence on proof generating processes that are neither structurally nor algorithmically transparent.

## 3.1 Opacity of Deep Learning

It is often said that DNNs lack epistemic transparency. It is important, however, to distinguish the *training* of a DNN from fully *trained* models. The procedure for training a DNN is algorithmically and structurally transparent. In most simple cases, it is algorithmically transparent that training works through the minimization of various loss functions through the iterative updating of weights on the connections between parameters in the model by means of the backpropagation of error gradients. There are extensive, well-maintained repositories that contain various structurally transparent implementations for training a wide variety of network architectures, and students taking a first class in machine learning are often required to write their own implementations of, say, the backpropagation algorithm (which, itself, involves a structurally transparent implementation of the chain rule for calculating the derivative of composed functions). There is nothing mysterious or opaque about *how* users go about training such models.

---

[4]Creel's treatment of computational transparency can be seen as a refinement and extension of computational concepts of understanding going back to Marr [Mar10].

[5]In cases where a computational system is algorithmically and structurally transparent (as in the proof of the Four Color Theorem), the reliability of the computational system at run-time can be (defeasibly) trusted [FR09, Due22] (though, in practice, one cannot transparently inspect it[Hum09]).

However, fully *trained* DNNs are said to be epistemically opaque [Hum04, Bog22], meaning that all of the epistemically relevant factors governing the model's behavior are fundamentally unsurveyable. In general, it is not possible to determine or 'fathom' [Zer22] in any intelligible or meaningful sense the algorithmic rules or principles governing the transformation of inputs to outputs of the model. As such, DNNs are opaque at the algorithmic level. This lack of transparency is due to the extremely high dimensionality and nonlinearity of the model, as well as the autonomous, error-driven, and semi-stochastic processes of weight assignment in guiding the settlement of the final parameterization of model.

Of course, with some DNNs we can say at any moment which graph is being analyzed, and there is a numerically trivial sense in which the trained model is fully transparent insofar as the values on the weights themselves are available to inspection (though not surveyable)[Lip18, Due23]. However unlike Appel and Haken's program, we cannot say *how* that graph is being evaluated, and thus the DNN is algorithmically and structural opaque.

## 3.2 Mathematical Knowledge with DNNs

Suppose that, as in the case of the Four Color Theorem, a mathematician wanted to know whether every graph in a set of graphs is reducible. Approaching this problem with a deep neural network, the mathematician trains a model on a large set of graphs known to be reducible and a large set of graphs known not to be reducible. The model is then evaluated (in the usual way) on a collection of graphs not included in the training set. It is shown to correctly classify every graph on which it is evaluated as reducible or not reducible. Suppose further that the model has never been shown to have misclassified the reducibility of a graph.

At this point, the mathematician unleashes their model on the graphs that constitute the possible counterexamples to the Four Color Theorem. After a minute or so, the model returns a result indicating that all them are reducible (and, thus, four-colorable). However, given that the model is not algorithmically transparent, it is also not mathematically transparent. Thus, we cannot get genuine mathematical knowledge from the model directly, as we fail to acquire an apriori warrant for the reasons we have discussed.

However, we do get strong, inductively justified belief in the Four-Color Theorem. Such a result bears some resemblance to a case [DVB+21] considered by Duede [Due23], where mathematicians use a DNN to guide mathematical attention to promising connections that led to the formulation and proof of a theorem linking specific algebraic and geometric properties of low-dimensional knots. Cases of this kind exemplify the potential for AI to assist mathematicians in their search for the most promising conjectures but leave the actual proof of the conjectures to humans. In this case our knowledge is not a direct result of the DNN, insofar as the ultimate responsibility for the proof lies with the human mathematician.

Consider, however, a hypothetical case in which a DNN trained to classify graph reducibility classifies all but one of Haken and Appel's 1,834 graphs as reducible, except for one which it classifies as irreducible. Here, the model has suggested a counterexample to the Four Color Theorem. Let us suppose that whether it is a genuine counterexample is something we can check, and that we check it, and find that it is, indeed, a counterexample. In this case too we now have genuine mathematical knowledge of a mathematical fact (namely, the falsehood of

the Four Color Theorem), even though again this knowledge cannot be said to follow *directly* from the output of the DNN, as it required human verification.

### 3.3 Mathematical Knowledge with LLMs

LLMs are particularly useful for mathematics as they output reports that are potentially linguistically and mathematically intelligible. However, like DNNs, LLMs are algorithmically and structurally opaque, and so they are afflicted by the epistemic limitations discussed in the previous section.

A recent case leveraging LLMs to achieve mathematical breakthroughs in extremal combinatorics involves a treatment of the Cap Set Problem in [RPBN+23]. A cap set is a subset of $(\mathbb{Z}/3\mathbb{Z})^n$ for which no three distinct elements sum to 0 (mod 3). For each $n$, the problem is to determine the size of the largest cap set. It is known that this number must be less than $\leq 3^n$ [Gro19], but its exact value is only known for $n \leq 6$. Moreover, the complexity of the solution space explodes for greater values of $n$ and so brute-force computational approaches are not feasible.

In [RPBN+23], researchers leverage a LLM to construct a cap set of size 512 for the case $n = 8$, a result that is significantly greater than the previously known largest value of 496. The approach begins by specifying an evaluation function that scores a candidate solution, where a solution is actually itself a Python program for generating a potential capset. The LLM then outputs a candidate Python program that is executed and scored by the evaluation function. If the program executes sufficiently quickly, it is sent to a program database. The system then samples the database and passes prior output programs to the LLM as inputs to repeat the generative process. This iterative approach generatively 'evolves' candidate programs. Eventually, this process identified a cap set of size 512, which human mathematicians verified to be correct.

Unlike in the Four Color Theorem, the solution generating procedure used here is not mathematically transparent. However, a human mathematician can easily survey and check its output. In this case, we get apriori knowledge that for $n = 8$ there is a cap set of size 512. Nevertheless, as this involves a human mathematician verifying this fact, we cannot say that genuine mathematical knowledge has been obtained directly from the computer.

### 3.4 Main Claim

With all these examples in mind, we claim the following.

<u>Main Claim 1</u>: If we want genuine mathematical knowledge directly from the testimony of a computer, then what the computer is doing must be mathematically transparent to us (as in the case of the Four Color Theorem). If what a computer is doing is *not* transparent to us (as in the case of typical DNNs or LLMs) then we cannot *directly* get genuine mathematical knowledge from the output of a computer, even though we can directly gain a type of inductively justified belief from it. However, even if we do not directly acquire mathematical knowledge from the testimony of a computer, if the computer outputs a human-checkable proof, example, or counterexample, then upon checking it we do gain genuine mathematical knowledge (though

not directly, insofar as human checking was required.)

## 4  Transparent Proof Checking

The considerations of the previous section might be taken to entail that, while extraordinarily useful, DNNs and LMMs can ultimately only be of circumscribed use in the generation of mathematical knowledge. In general, their reports only offer us inductively justified belief (e.g., a kind of empirical result), and it is only when they output mathematically transparent results that can be human-checked that we can acquire mathematical knowledge from them.

However, there is a fairly straightforward way to surpass these limits. Let us focus on the case in which a machine (perhaps an LLM) outputs not only some mathematical result $X$, but also outputs something resembling a proof of $X$, all of which is stored somewhere on a hard drive. If what has been stored on the drive can be human-checked, then we can check it, and if it is indeed a correct proof, we thereby gain apriori knowledge of $X$.

The limitation appears to arise if, as in the case of the Four Color Theorem, the proof is so complex that it is not human-checkable. As we argued in Section 3, if the result is not generated by a mathematically transparent process, the output of the computer only gives us inductively justified belief in the truth of $X$.

But, this does not mean that apriori knowledge of $X$ lies beyond our reach. Let us suppose that the stored proof of $X$ is such that, while enormously complex, it is systematically organized as a sequence or tree of propositions of the sort one might encounter in a mathematical logic class. We can imagine constraining the output of the machine generating the proof in such a way as to demand that its outputs be formulated in this way (as in [RPBN⁺23] where the model outputs all results in syntactically correct `Python`). We can generously allow various abbreviations, additional rules, and verbose articulations of the steps in the proof so that this proof has roughly the form of a human generated proof[6], even though it was not generated by a human and is in fact so large that it cannot be surveyed by a human. These assumptions can be made to hold by adding some overhead to the original program and forcing the proof of $X$ to be stored in this form.

While we cannot check this proof, this does not stop us from writing a proof-checking program that can. The proof-checking program we imagine simply goes through the proof, verifying that it starts with genuine axioms, and then verifying that each step is a legitimate application of some standard logical or mathematical rule. We can imagine a version of this proof-checker that is, in fact, completely mathematically transparent, approaching the problem of checking the proof in exactly the way a human would. When such a program runs we can correctly say at any point something like 'the computer is now checking step 154835 , and is verifying that it is a correct application of modus ponens'.

Let us assume that we run the proof-checker and it reports no errors. Just as Appel and Haken can acquire an apriori warrant for believing the Four Color Theorem from the output of their mathematically transparent program, so too we can acquire an apriori warrant for

---

[6]Of course, the proofs of practicing mathematicians are not like this, often involving large leaps and a kind of hand wavy skipping over of 'trivial' steps [Kit98].

believing that there is a correct proof of $X$ from the output of our mathematically transparent proof-checker. From the fact that there is a correct proof of $X$, the truth of $X$ follows, and thus we have acquired an apriori warrant for believing $X$. If everything has gone well we can even come to *know* that $X$ is true. This outcome is true even though no human has (or ever could have) any sort of rational grasp on the process that led to the generation of the proof, and no human is capable of checking the proof stored on the hard drive.

Of course, if the LLM 'claims' to have proven $X$ but cannot produce and store the actual proof of $X$, then we cannot use a proof-checker in the way just described. In this case, we see no way to acquire anything other than inductive grounds for believing $X$.

This leads us to the second of the two central claims of this paper, which may be viewed as a counterpoint to Main Claim 1.

<u>Main Claim 2</u>: We *can* (indirectly) gain apriori knowledge from the output of a computer program that is not mathematically transparent but which stores a (not necessarily human-checkable) proof of a mathematical claim. This is accomplished by employing a mathematically transparent proof-checker to evaluate the stored proof of the claim.

## 5   General Conclusions

Modern LLMs and DNNs are opaque to us in ways that create obstacles to obtaining mathematical knowledge from them. However, we have argued that if a proof-checker transparently automating human forms of mathematical evaluation is attached to such machines, then we can obtain apriori mathematical knowledge from them. Surprisingly, this applies even in cases where the original machines are entirely opaque to us and the proofs they output are not human-surveyable.

A question for further consideration is to what extent we may gain scientific [BKLW23] knowledge outside of mathematics by appending analogous transparent 'checking' mechanisms to the output of otherwise opaque algorithms. This would get us closer to overcoming the perceived problems of confabulation, and realizing the ambition of fully automated scientific discovery.

## References

[AH89]   Kenneth I Appel and Wolfgang Haken. *Every planar map is four colorable*, volume 98. American Mathematical Soc., 1989.

[BKLW23]   Abeba Birhane, Atoosa Kasirzadeh, David Leslie, and Sandra Wachter. Science in the age of large language models. *Nature Reviews Physics*, 5(5):277–280, 2023.

[Bog22]   Florian J Boge. Two dimensions of opacity and the deep learning predicament. *Minds and Machines*, 32(1):43–75, 2022.

[Bur98]   Tyler Burge. Computer proof, apriori knowledge, and other minds: The sixth philosophical perspectives lecture. *Philosophical perspectives*, 12:1–37, 1998.

[Cre20]     Kathleen A Creel. Transparency in complex computational systems. *Philosophy of Science*, 87(4):568–589, 2020.

[DL80]      Michael Detlefsen and Mark Luker. The four-color theorem and mathematical proof. *The Journal of Philosophy*, 77(12):803–820, 1980.

[Due22]     Eamon Duede. Instruments, agents, and artificial intelligence: novel epistemic categories of reliability. *Synthese*, 200(6):491, 2022.

[Due23]     Eamon Duede. Deep learning opacity in scientific discovery. *Philosophy of Science*, 90(5):1089–1099, 2023.

[DVB⁺21]    Alex Davies, Petar Veličković, Lars Buesing, Sam Blackwell, Daniel Zheng, Nenad Tomašev, Richard Tanburn, Peter Battaglia, Charles Blundell, András Juhász, et al. Advancing mathematics by guiding human intuition with ai. *Nature*, 600(7887):70–74, 2021.

[FR09]      Roman Frigg and Julian Reiss. The philosophy of simulation: hot new issues or same old stew? *Synthese*, 169(3):593–613, 2009.

[Gro19]     Joshua A. Grochow. New applications of the polynomial method: the cap set conjecture and beyond. *Bull. Amer. Math. Soc.*, 56(1):29–64, 2019.

[Hum04]     Paul Humphreys. *Extending ourselves: Computational science, empiricism, and scientific method*. Oxford University Press, 2004.

[Hum09]     Paul Humphreys. The philosophical novelty of computer simulation methods. *Synthese*, 169(3):615–626, 2009.

[Kit98]     Philip Kitcher. Mathematical change and scientific change. *New directions in the philosophy of mathematics*, pages 215–242, 1998.

[Lip18]     Zachary C Lipton. The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery. *Queue*, 16(3):31–57, 2018.

[Mar10]     David Marr. *Vision: A computational investigation into the human representation and processing of visual information*. MIT press, 2010.

[RPBN⁺23]   Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Balog, M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang, Omar Fawzi, et al. Mathematical discoveries from program search with large language models. *Nature*, pages 1–3, 2023.

[TWL⁺24]    Trieu H Trinh, Yuhuai Wu, Quoc V Le, He He, and Thang Luong. Solving olympiad geometry without human demonstrations. *Nature*, 625(7995):476–482, 2024.

[Tym79]     Thomas Tymoczko. The four-color problem and its philosophical significance. *The journal of philosophy*, 76(2):57–83, 1979.

[Tym98]    Thomas Tymoczko. *New directions in the philosophy of mathematics: An anthology*. Princeton University Press, 1998.

[Wil13]    Timothy Williamson. How deep is the distinction between a priori and a posteriori knowledge?   In Albert Casullo and Joshua C. Thurow, editors, *The A Priori in Philosophy*, pages 291–312. Oxford University Press, 2013.

[Zer22]    John Zerilli. Explaining machine learning decisions. *Philosophy of Science*, 89(1):1–19, 2022.