

On the relevance of logic for AI: misunderstandings in social media, and the promise of neuro-symbolic learning

Vaishak BELLE^a

^a *University of Edinburgh & Alan Turing Institute, UK*

Abstract In this position paper, we examine some of the assumptions held and circulated in social media about logic and its relevance to the development of modern AI, which is primarily driven by deep learning. The paper aims to address fundamental misunderstandings about logic and ultimately argue for the benefits of symbolic formalisms in modeling uncertain worlds. While it is now recognized that statistical associations learned from data are limited in their ability to understand the world, there is still a great deal of criticism and hesitancy regarding the use of symbolic logic to achieve or support a broader vision for AI. By arguing that symbolic logic is more flexible than non-experts believe, we make a case for Neuro-Symbolic AI, which offers the best of both worlds.

Keywords. neuro-symbolic AI
logical foundations
GOFAI

Preface

This article lies between a position paper and a survey paper, and it does two things. On one hand, it discusses the breadth and diversity of solutions encompassed within symbolic logic. We believe that many of these dimensions are not obvious to people outside the logical community, and perhaps, even those working within certain areas of logic might not be aware of the latest developments in statistical relational learning. On the other hand, it points to common objections to using logic when building complex AI systems involving machine learning. This reflects objections raised by Geoff Hinton and echoed by others. Essentially, this paper serves as a survey that tackles both of these aspects.

1. Introduction

Artificial Intelligence (AI) is widely acknowledged as a new kind of science that will bring about (and is already enabling) the next technological revolution. Virtually every week, exciting reports come our way about the use of AI for drug discovery, game playing, stock trading and law enforcement. And virtually all of these are mostly concerned

with a very narrow technological capability, that of predicting future instances based on past instances.

Identifying statistical patterns, correlations, and associations is, without doubt, extremely useful. In the first instance, it is needed in numerous applications to inspect features and properties of interest in observed data. It serves as the backbone of recommendation systems, for example, and is likely more than sufficient, even with flaws, when gathering context. While searching for “how to raise lambs” in an online bookstore, we might be a little disappointed if it suggests “Silence of the Lambs” by Thomas Harris, and somewhat annoyed if it suggests cookbooks on “how to cook lamb”, but such low quality results are unlikely to have long-term effects. This type of AI might also be useful but somewhat problematic for, say, fast-tracking the review of job applications, provided these models are adjusted for bias, and a human intervenes and interprets the outcome and determines how to act further. This type of AI was largely believed to be sufficient for vision systems [172], until it was observed that self-driving cars fail stupendously, and that the state-of-the-art systems can be fooled in strange and unnatural ways [75].

Be that as it may, this is a very narrow view of AI capabilities. AI, as understood by both scientists and science fiction writers, is clearly much broader. What distinguishes big-data analysis from AI is that the set of capabilities we wish to enable with the latter. We are interested in a general-purpose, autonomous computational entity that, in the very least, has agency. Many of these concerns were widely debated, discussed, and developed during the heyday of good old-fashioned AI [51,114,110].

However, despite recognizing that data-driven statistical learning is limited in its ability to understand the world and model its knowledge [123], there is still a lot of criticism and hesitancy about the use of symbolic logic to accomplish or assist in a broader vision for AI [44].

In this position paper, we examine some of the assumptions held and circulated in social media about logic and its relevance to the development of modern AI, which is primarily driven by deep learning. The paper aims to address fundamental misunderstandings about logic and ultimately argue for the benefits of symbolic formalisms in modeling uncertain worlds. By arguing that symbolic logic is more flexible than non-experts and critics believe, we make a case for Neuro-Symbolic AI, which offers the best of both worlds.

2. Why Tweets?

Before going further, let us briefly reflect on the objects of focus – mainly, tweets from Twitter – in this paper.¹

Our focus on views of AI from Twitter may seem unusual; however, Twitter has turned into a dominant space for public statements by leading experts.

One could argue that these do not constitute well-developed and scientific arguments. However, opinions from leading deep learning experts about logic is unlikely to appear in any peer-reviewed venue as logic is largely dismissed by many from the main-stream machine learning community, such as Geoff Hinton. In fact, one could argue that

¹Twitter was recently rebranded as X, but we will continue to refer to the social media website as Twitter for ease of readability as we often use capital letters as variables in logic. Posts on Twitter, as usual, are referred to as Tweets.

precisely because these opinions are not based on scientific arguments, they are held implicitly and cannot be easily challenged, except in the manner here, or yet again, informally on social media websites. So, although such views are not as iron-clad as peer-reviewed position papers, we view them as scientific opinions all the same: positions expressed by peers for peers.

Of course, it is possible some may wish to retract statements made in tweets, saying that the space limitations forced them to make informal remarks that could be easily misunderstood, or that there was implicit irony. This is why we have included screenshots in most instances, or otherwise linked to them, and admit that we are taking those statements at face value and apologise if we have misrepresented an individual's position.²

The downside of including only a few tweets is that they cannot accurately represent the entire community due to the small sample size. We acknowledge that this limitation cannot be avoided without comprehensive surveys involving a sizable portion of the community. For example, a recent survey in [77] analyzed expert opinions on the technological progress of AI surpassing human performance. We believe that such surveys could be a possibility for future work. However, for these surveys to be effective, we need to clearly define the positions and technical aspects we are seeking expertise on from the experts. By examining some notable tweets from leading experts, we aim to highlight key points that are often misunderstood, which could then be incorporated into such surveys.

A more effective way to rebut misunderstandings is to provide concrete technical demonstrations that contradict claims, especially negative claims made by critics. We believe these types of demonstrations are already being discussed in the Neuro-Symbolic AI community. For example, the work in [184] characterizes learnability results for a popular weak supervision demonstration, e.g., as seen in DeepProbLog [121]. The idea is that instead of directly labeling MNIST images, we assume that labels are provided for a logical or arithmetical operation on the numbers represented by these digits. The result of this operation is then used as feedback for the neural network, which is then evaluated on its ability to correctly interpret the samples. Such a pipeline is difficult to implement in general without a symbolic reasoner. Likewise, in [69,87], it is shown that logic-based loss function approaches not only guarantee that deep learning predictions satisfy the constraint but also achieve this satisfaction with far fewer samples, making them more efficient.

However, existing approaches and demonstrations cannot fully capture the wide-ranging criticisms of logic found on social media. By its nature, the scope of a scientific paper is limited in order to be rigorous. This is why we discuss a variety of negative claims and misconceptions in this paper. While it may be impossible to be exhaustive, we believe we have addressed most of the dimensions in which logic is commonly viewed as constraining.

What is particularly notable is a considerable overlap in ideas and technical constructs across sections in this article. This is because there are a number of intriguing connections between logic, probability and learning. Readers interested in this should refer to the following surveys: a book-length treatment on statistical relational learning in [149], the interplay between logic and probability in [177] and [16], between logic and learning in [17], and between logic and deep learning in [18].

²We have freely used screenshots of tweets without reaching out to the users for permission, assuming that these tweets can be cited like a website since they are public and have an openly accessible URL.



Figure 1. Hinton’s analogy

3. Logic is old-fashioned

In the first part of this article, we will look at some of the criticisms against using logic. We then turn to a number of positive dimensions to examine the integration of logic and learning.

3.1. *Neural approaches and nothing else!*

Modern AI has moved on, we are told. The idea of using symbolic logic is outdated, and the area of knowledge representation defined over symbolic logic is now affectionately (or perhaps pejoratively) called good old-fashioned AI, or GOF AI for short.

In the early days of AI, John McCarthy put forward a profound idea to realise artificial intelligence (AI) systems [124]: he posited that what the system needs to know could be represented in a formal language, and a general-purpose algorithm would then conclude the necessary actions needed to solve the problem at hand. The main advantage is that the representation can be scrutinised and understood by external observers, and the system’s behaviour could be improved by making statements to it.

Numerous such languages emerged in the years to follow, but first-order logic remained at the forefront as a general and powerful option [132]. Propositional and first-order logic continue to serve as the underlying language for several areas in AI, including constraint satisfaction [28], automated planning [155], database theory [116], ontology specification [103], verification [11], and knowledge representation [114].

And yet, “modern” AI has decided that these efforts are superfluous, or at least easily replaceable once a training dataset has been created. As an infamous and inflammatory instance, Turing-award winner Geoff Hinton remarked that fixating on symbols was a waste of time, analogous to funding research on gasoline engines. Implicit here is the argument that we clearly need to be focussing on electric engines, presumably analogous to deep learning.³

In 2020, he reiterated his position and suggested that:⁴ “Deep learning is going to be able to do everything.” Strangely, his position seems to have changed over the years,⁵ but

³<https://twitter.com/tabithagold/status/1070736319901519876>

⁴<https://www.technologyreview.com/2020/11/03/1011616/ai-godfather-geoffrey-hinton-deep-learning-will-do-everything/>

⁵<https://www.noemamag.com/deep-learning-alone-isnt-getting-us-to-human-like-ai/>

it is hard to get a sense of what kind of mixture of symbols vs learning he is advocating for. For example, in a very recent interview after quitting his position at Google, the following transpired:⁶

The dominant idea at the time, known as symbolic AI, was that intelligence involved processing symbols, such as words or numbers.

But Hinton wasn't convinced. He worked on neural networks, software abstractions of brains in which neurons and the connections between them are represented by code. By changing how those neurons are connected — changing the numbers used to represent them — the neural network can be rewired on the fly. In other words, it can be made to learn.

“My father was a biologist, so I was thinking in biological terms,” says Hinton. “And symbolic reasoning is clearly not at the core of biological intelligence.”

“Crows can solve puzzles, and they don't have language. They're not doing it by storing strings of symbols and manipulating them. They're doing it by changing the strengths of connections between neurons in their brain. And so it has to be possible to learn complicated things by changing the strengths of connections in an artificial neural network.”

The “theory of everything” approach in science, or perhaps its analog in AI, that of having a single algorithm/architecture/framework for all tasks [55], is undoubtedly appealing. Some theoretical physicists have hopes pinned on string theory, for example, to come up with a single framework that unifies all observational data, across large and minuscule physical bodies [53]. Likewise, the appeal of purely neural model is attractive. However, there is lots to debate here.

Firstly, deep learning models are loosely inspired by the brains but not fully accurate representations (yet) [168,131]. Secondly, there is the notion of innateness [173], and how much evolution might help the brain in understanding and processing the world in a structured manner. And thirdly, we must bear in mind that we still lack a complete understanding of how the neurons of a bird (let alone a human) are wired, and how that influences cognitive capabilities. Merely knowing that neural weights enable birds to solve puzzles and recognize faces does not necessarily imply that our implementation of their neurons should resemble or possess similar properties. These concerns had also been debated in the literature in the 1980s [152,165].

Putting such issues aside, it is also worth noting that proponents of the symbolic approach to AI never explicitly claimed the existence of symbolic representations within our minds [30,114]. In essence, the symbolic approach offers a coherent strategy for: (a) executing symbolic expressions, which capture the knowledge of the system about the world, and (b) comprehending the (idealized) implications of one's knowledge, as specified by inference rules in logic.

As argued by Levesque [113], this is not a novel concept – Leibniz articulated centuries ago that certain types of thinking adhere to symbolic processing. Hence, why not employ an algebraic treatment for cognition? As scientists, we may debate whether it is more useful to have an exact model of computation that approximates the reasoning in the brain [165,95] or whether we should forego these models altogether and simply be satisfied with informal descriptions of reasoning [146], as might emerge from a trained model [43].

⁶<https://www.technologyreview.com/2023/05/02/1072528/geoffrey-hinton-google-why-scared-ai/>

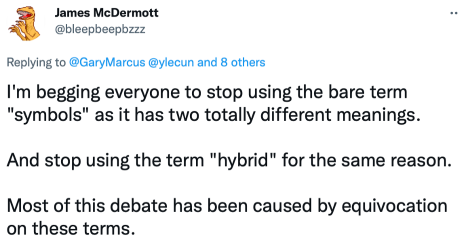


Figure 2. DcMermott on symbols.

We reiterate that the allure of a purely neural approach is understandable, given its simplicity and the sense of a “unified theory” it evokes. However, the arguments regarding the effectiveness of the training process in capturing intricate reasoning [87] and the potential for incorrect [176] and unreliable predictions [5] suggest that a purely neural approach may not be sufficiently robust to exploit and capture structure.

By taking a step back, we realize that until the past few centuries, our understanding of the brain and neurons was limited. Yet, during this time, we were able to calculate, develop number theory, construct calculators, and ultimately build computers [174]. Imagine if we had solely dedicated ourselves to constructing elaborate brain replicas in the hopes that they could handle (say) tax calculations for us. Most importantly, we cannot test for a capability without first defining that capability, such as (say) deduction [146].

All of this underscores the significance of the symbolic approach, which offers an idealized framework for well-defined (relative to the formal language) forms of reasoning. There is a popular analogy [30] suggesting that we need not build wings and feathers to build airplanes; comprehending the principles of aerodynamics is enough. So, why shouldn’t the development of a theory of artificial cognition be just as relevant for a type of AI that is behaviorally similar to humans in some instances, without necessarily resorting to a brain-like architecture?

Geoff Hinton, of course, is not alone in being dismissive about symbols. There are many other severe views on the relevance of logic for modern AI, and in what follows, we will survey and respond to a selected set of misunderstandings extracted from Twitter.

3.2. *There is a dichotomy*

A common view held by many in the broader community that there is an inherent dichotomy between symbolic logic and machine learning, the former focused on discrete structures and the latter focused on continuous representations. In fact, even scientists within the AI community make this distinction [154], and suggest that logic is not really appropriate for machine learning. Consider, for example, the tweets by James McDermott.

This, admittedly, is not even necessarily negative on the topic of symbols, but just points out that: (a) symbolic logic as used in AI is focused on discrete symbols; and (b) symbolic processing in vector (real-valued) space is a separate topic of study that can be independently done from symbolic logic.

What we are seeing here is a narrowing of the use of “logic” simply as *classical logic* – say, as introduced in [59] – defined over Boolean truth values. Moreover, the

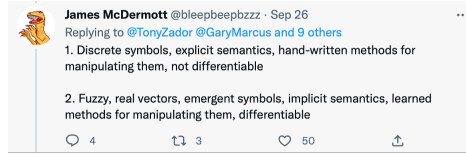


Figure 3. DcMermott on dichotomies.

use of logic is also assumed to be limited to discrete propositional assertions, as seen in ontologies that capture relationships and hierarchies about commonsensical concepts [125], as well as in early attempts at logic programming [104].

We will now discuss the use of non-Boolean truth values and continuous properties in logic, and how that is making an appearance in the area of neuro-symbolic AI.

3.2.1. Real-valued truth values

To a large extent, it is true that the area of knowledge representation in AI focuses on discrete symbols and a Boolean interpretation [30]. But, on the other hand, it's been close to 60 years since we have fuzzy logic [190], among others languages for non-binary truth values [100]. These allow us to assign a truth value between 0 and 1 to propositions, with the understanding that these values indicate the degree to which the proposition may be true. Fuzzy logic can also be utilized to represent ambiguous concepts, such as stating that a person is tall, without specifying tall as a categorical property.

The use of such values to propositions means that the interpretation of Boolean connectives also changes. For example, the formula $\alpha \wedge \beta$ could be mapped onto the *min* of the truth values of the individual formulas. That is:

$$\alpha \wedge \beta \doteq \min(\alpha, \beta).$$

If α is assigned a truth value of 0.6 and β is assigned a truth value of 0.4, then the conjunction would be given a truth value of 0.4. Of course, one can see that if the truth values are either 0 or 1, then the *min* function aligns with classical logic in the sense that if either α or β is 0, then the conjunction will also have the truth value of 0.

By construction, the outputs of neural networks can be mapped to real numbers between 0 and 1. Owing to the nature of truth values in such logics, these outputs can be directly modeled as atoms in logical formulas. This led to an early wave of *neurosymbolic AI* formalisms [70] and the development of a field that integrates neural outputs in a logical language [86]. Perhaps the most representative examples in this space are logic tensor networks [8] and other approaches based on fuzzy logic [180]. The motivation for many of these languages is to logically capture concepts that have been learned from neural networks, in order to reason about these concepts as part of a commonsensical knowledge base. Thus, the agent would be reasoning about hierarchies and relationships, but many of the relations in this knowledge are learned directly using neural networks, presumably from observational data.

It is worth noting that reasoning about concepts and relations is an ongoing problem with neural networks – see efforts such as capsule networks [158] and module networks [4] – and there are very few general solutions. Neuro-symbolic AI is stepping in here,

especially if it were to allow a general framework for injecting knowledge expressed in a fragment of first-order logic, could be very welcome.

3.2.2. *From discrete to continuous*

Capturing the output of neural networks as truth-values in a logical formula is one approach to reasoning about vector spaces. However, we can also use logic to reason about continuous properties as formulas.

Although it is common to discuss discrete properties in logical AI, it is not necessary that they must do so. Logical formulas are indeed discrete structures, but they can also express properties about countably infinite or even uncountably many objects [84,151,25,17].

Reasoning about real numbers have long been an area of interest in mathematical logic [96], going back to Tarski, and are a major concern in satisfiability modulo theories (SMT) [11]. SMT can be seen as a generalization of SAT for propositional logic and is being used for the verification of timed and hybrid systems that involve both discrete and continuous properties. For example, the following formula expresses that a logical function with one argument f applied to x is lesser than the square of y . This could be conjoined using Boolean connectives with other assertions, such as one that says y is greater than the two-variable function g applied to x and z :

$$f(x) \leq y^2 \wedge y > g(x, z).$$

Here, the domain of x, y , and z could range over the set of natural numbers \mathbb{N} , the set of integers \mathbb{Z} , or even the set of reals \mathbb{R} .

Therefore, we can use these formulas to represent constraints on geometric spaces. A recent body of work has examined the idea of regularizing neural networks by adding logical constraints to the loss functions. The idea is to train the network such that the loss is calculated against this logical constraint, which is backpropagated. The goal then is to train the network in such a way that predictions always satisfy these logical constraints. There is existing work on propositional constraints [69], real-valued constraints [87] as well as temporal formulas [93], the latter of which trains the network to dynamically navigate an environment in only the valid geometric space.

One of the interesting observations in almost all of these papers on loss functions is that they demonstrate that it is much more effective to train the network using such loss functions than assuming the constraints are represented in the data. So it is much more sample efficient [92]. Moreover, some of these architectures also allow for the complete satisfaction of the constraints [87]. This is necessary in safety-critical and high-stakes applications.

3.3. *Logic is not good for probabilistic uncertainty*

Classical quantifiers in logic, as well as the connectives, allow for disjunctive uncertainty, the existence of individuals, and properties applicable to all individuals in the domain. Because the data we collect is often noisy, or we sometimes have to approximate and average over populations, the use of probability theory is essential [141]. Since classical logic traditionally did not represent probabilistic assertions, much of the learning and

uncertainty in the AI community moved away from logic. We will discuss here that the connection between logic and probability is deep, and there is a vibrant community focussed precisely on this agenda [149].

3.3.1. Probabilistic logical models

Since the work of Nilsson [134], the use of logic to capture non-trivial probabilistic spaces and reason logically about events in those spaces has been a major concern in uncertainty quantification in AI [157] and statistical relational learning [149]. The key idea here is that it should be possible to assign probabilities to atoms, which would then provide a way to extend these probabilities to complex formulas. That is, if α is a well-defined (classical) formula in a logical language \mathcal{L} , then so is $\text{Pr}(\alpha)$ [81]. This leads to a representation language that may involve a combination of deterministic and probabilistic assertions, capturing the knowledge base of a putative agent. For example [24], consider the following formula:

$$\alpha \wedge \text{Pr}(\beta) > \text{Pr}(\gamma) \wedge \text{Pr}(\gamma) \leq 0.6.$$

It is assumed that α is true, and the probability that β is true is greater than the probability that γ is true. Additionally, γ is believed with a probability of less than or equal to 0.6. Here, α may be a non-probabilistic assertion. The probability of γ is not given a unique value, and we are allowed to compare the likelihoods of two formulas. Such combinations are difficult to express using probability theory alone.

In recent years, there has been a steady progress on designing languages that can not only capture Bayesian networks and factor graphs [107], but also extend them with a relational and a logical syntax. Popular languages for pragmatic specifications of logic and probability include Markov logic networks [154], ProbLog [150] and BLOG [129]. Many of these not only investigate the representational restrictions that enable the capture of distributions succinctly, but also explore how to reason with the resulting distribution, and in some cases, learn the distributions or representations themselves. (They have to restrict the expressiveness of the language in order to ensure that their representations capture a single distribution; so the above formula may be difficult to express here too.) Consider the following program in ProbLog [150]:

```
0.5::heads1.  
0.6::heads2.  
twoHeads :- heads1, heads2.
```

This allows us to capture a mixture distribution composed of a biased coin toss and an unbiased coin toss, with the latter having a 0.6 probability of landing heads.

Interestingly, Bayesian networks can also be modeled as ProbLog programs [149]. And what is more interesting is that probabilistic inference in Bayesian networks [38], ProbLog programs [65], Markov logic networks [154], and factor graphs [107] can all be shown to be reducible to the same computational task known as *weighted model counting* [6]. Weighted model counting is an extension to SAT in the sense that each satisfying assignment is assigned a weight. By computing the sum of the weights of all satisfying assignments, we can relate that sum to the conditional probability and marginals in a

Bayesian network. That is, for a propositional language \mathcal{L} , assume a weight function w maps its literals to $\mathbb{R}_{[0,1]}$. Then, for some $\phi \in \mathcal{L}$,

$$\text{WMC}(\phi, w) = \sum_{\{M \mid M \models \phi\}} \prod_{\{l \mid l \in M\}} w(l).$$

The product operation here is defined in terms of all the literals that are true in a given model of ϕ .

As argued in [178, 16], it is not only the case that logical languages allow us to reason about probability distributions over combinatorial spaces, but it is also the case that the syntax of logic can help capture complex relationships that are difficult to model using standard probabilistic languages [73]. Moreover, by way of weighted model counting, there is a single generic approach for probabilistic reasoning over discrete, combinatorial spaces that is competitive [38]. It is also amenable to both exact as well as approximate inference schemes [36].

Recently, there have also been extensions from discrete combinatorial spaces to continuous ones [26, 40], referred to as *weighted model integration*. Here, the formula $x \in [-5, 5]$ with a weight of 0.56 might represent a continuous random variable x whose piecewise constant density for all values between -5 and 5 is 0.56. Analogously, the same formula with the weight of $x^2/2$ might represent a piecewise polynomial density specification for x , such that for all values between -5 and 5 , its density is given by the square of that value divided by 2. As with weighted model counting, inference in this formulation is performed by means of a notion of model counting in SMT [11].

3.3.2. Generalising the specification of a distribution

Going back to the history of the use of logic in AI [132], there has been considerable interest in unifying logic and uncertainty. Note that, through the use of quantifiers, it is possible to express uncertainty that may not always align with a single distribution. For instance, McCarthy [126] was concerned about probabilities in the early years of using first-order logic for knowledge representation. However, he makes a very salient point that we need to think carefully how numbers and first-order sentences fit together. For example, he argues [126]:

(i) *It is not clear how to attach probabilities to statements containing quantifiers in a way that corresponds to the amount of conviction people have.*

(ii) *The information necessary to assign numerical probabilities is not ordinarily available. Therefore, a formalism that required numerical probabilities would be epistemologically inadequate.*

His point, simply, is that we should not be expected to put probabilities on every formula; sometimes it suffices to say that $p \vee q$ holds without saying which, and by how much. Moreover, if we assign a probability of r on that formula, or to, say, $\exists x P(x)$, such an assertion in itself does not provide any additional information on how to further assign a probability to p , q , $P(a)$, and so on. Many popular languages for logic and probability mentioned above, including Markov logic networks [154], ProbLog [150] and BLOG [129], do not allow this level of flexibility. In fact, this requires a different type of machinery altogether, one which permits multiple prior distributions [23]. Consider a sub-formula from the example from above:

$$\Pr(\gamma) \leq 0.6.$$

The formula should, in principle, allow for every distribution that accommodates a probability of γ being less than or equal to 0.6. In contrast, in ProbLog, it is assumed that there is a single distribution over the model, and not specifying a probability on (say) a disjunction might be interpreted as a hard constraint that is true in all possible worlds. However, there are languages that do permit such rich specifications. See, for example, works such as [137] and [24].

More generally, probability measures [68] on first-order structures and other proposals on logic and uncertainty [154,150,129,23] allow us to append probabilities and weights in a logical language in different ways, yielding formal frameworks that go beyond and generalize the standard definition for a probability space. There are also approaches [56] that are based on possibility theory, which permits a different model for uncertainty that can be powerful when experts disagree or are uncertain about probabilistic assertions.

3.4. *Symbols do not always need a logic*

In the machine learning literature, it is not uncommon to find syntactical objects, especially well-defined symbolic expressions, such as programs, that are learned without an explicit definition of the semantics [109]. In such cases, one would need to define only the interpreter and the compiler [58], with an implicit notion that the atomic objects refer to concrete objects in the real world, as obtained by the process of *symbol grounding* [170].

However, with programs in the program induction literature [79], there is (or rather, should be) an implicit logical syntax and semantics that defines: (a) what sort of expressions can be constructed, and (b) what they mean and capture. For example, sequential instructions could be understood as conjunctions, and while loops can be captured using second-order quantification [112,171,79]. If we further want to understand what properties are entailed by these programs, then we need to define the semantics comprehensively and analyze what follows from the logical theory corresponding to a program.

Indeed, without a clear specification of how compositions of expressions should be interpreted and evaluated, how are we to know what these programs are yielding [124]? There has been a surge of a new family of programming languages that capture intricate machine learning models. Typically, these languages allow the use of random primitives as well as operators for conditioning and providing evidence. These are referred to as *probabilistic programming* – see, for example, Church [76], ProbLog [49], and the generic construction in [166]. In some cases, they might support combinations of discrete and continuous distributions, and higher-order functions. A general approach to understanding how these programs can be constructed and what sort of distributions they model is through the use of a formal semantical setup, usually in a fragment of first or second-order logic.

See also works such as [12] for discussions on attempting to construct the semantics for one programming language syntax from another. Such a move is especially desirable if we want to check for the internal consistency of an ad hoc programming language. For philosophical arguments on the importance of semantics, see, for example, [42].

3.5. Logic is about categorical propositional assertions

As discussed above, often “logic” is synonymous with (the classical interpretation of) propositional logic.

There are many systems for writing down symbols, and interpreting logical symbols and formulas built up these symbols. Classical approaches include propositional logic (Boolean symbols, A and B is true iff A is true and B is true) and first-order logic, which uses quantifiers. In first-order logic, there is a domain of discourse which stands for the objects in the world. We then say that $\exists x. P(x)$ is true if and only if there is some individual from the domain of discourse such that the property P is true for that individual. Likewise, the formula $\forall x, y. Grandparent(x, y) \supset \exists z. Parent(x, z)$ says that if x is a grandparent of y , there must be some individual z whose parent is x . First-order logic can also use functions over reals, as seen in satisfiability modulo theory (SMT) [11]. As we illustrated with weighted model integration, which is also defined using SMT, we can express formulas such as $x \leq 5 \wedge x \geq -5 \wedge x > y^2$. Here, both the variables are assumed to be nullary functions. But we could also have functions with arguments and nestings of these functions to construct well-defined formulas of the sort: $f(f(x, y), y) \leq y^2$.

We might also be interested in entertaining multiple possible truth assignments to model uncertainty about the environment. For example, there is modal logic [106], which can capture possibilities, beliefs, and intentions [160]. A variant of modal logic with numbers on worlds can lead to probabilistic logics [81], that allow us to reason about probabilities on formulas [63] as well as beliefs about these formulas [62,23].

Beyond these formalisms that map atoms (and by extension, formulas) to binary truth values, there are logics that relax that assumption. Fuzzy logics map Boolean symbols to real numbers, leading to real-valued semantics for non-atomic formulas constructed using connectives. For example, if A and B get values between 0 and 1, then $A \vee B$ gets a value of 1 iff $\min(A, B)$ is 1. Moreover, the conjunction could also get a value between 0 and 1, by way of $\max(A, B)$. Such a definition reduces to the classical semantics when both A and B are assigned 1, in which case the maximum of the two would also be 1.

These are all part and parcel of symbolic logic. The choice of the language, the choice of the semantic rules that we use over the well-defined formulas, along with its computational properties such as decidability are aspects of a logical framework. Moreover, once a logical framework is considered, we could choose to prove logical entailments either by considering assignments to the variables and seeing if the consequent follows, or by applying inference rules established in a proof theory [83]. If we choose to add weights [38], measures [81], or belief functions [56], this then leads to notions such as weighted model counting [6] and algebraic model counting [101], defined over the models of a formula (i.e., possible worlds). Ultimately, we could consider theorem proving [83], model checking [9], SAT solving [11], or model counting [74], depending on the context and application.

Each of these dimensions is already impacting current inquiries into the properties of machine learning models. For example, to tasks from knowledge-based completion to reasoning with ontology triples using neural techniques, there has been development on so-called *neural theorem provers* [130]. These are inspired by Prolog’s proof-theoretic backward chaining mechanism [49] and the aim in those works is to implement that



Figure 4. Parikh on nonmonotonic logics.

scheme in an end-to-end learning paradigm. Both SAT solving [185] and model counting [69] are important ingredients in state-of-the-art approaches to regularizing neural networks using logical formulas. This is motivated by the need to ensure neural network predictions always satisfy certain domain constraints. Model checking tools are mainstream for checking the robustness of neural networks [78]. There is also some work [180] on studying whether using real-valued fuzzy logics to permit differentiability in neural networks is comparable to differentiability as a result of probabilistic extensions to model counting [69].

In summary, we can explore a variety of logical syntax and semantics, each of which may have interesting interactions with machine learning properties and capabilities.

3.6. Monotonicity

Classical logic is monotonic. That is, if $\alpha_1, \dots, \alpha_k \models \beta$, then it cannot be the case that adding new knowledge, say, α' forces us to retract β : formally, it has to be that $\alpha_1, \dots, \alpha_k, \alpha' \models \beta$ also.

John McCarthy was concerned about the problem of monotonicity and wondered how we might deal with exceptions and abnormality. The problem of monotonicity is so ubiquitous, it even comes up in the formulation of automated planning [153]. For example, imagine that you have an action to paint a box blue and another action that pushes the object. Let us say we paint the object and next, we push the object. When we execute the second action, it is implicit that the color of the object does not change. So we would have to somehow codify not only what the effects of the push action are, but also what the non-effects are. And if we start writing down all the non-effects, there could be exponentially many. Moreover, there are various preconditions that must hold for us to be able to push the object. For instance, we should be strong enough to push it, we must not be holding other objects, we are presumably operating under reasonable gravity assumptions, and so on. And if we start expressing all of them, it again looks like a hopeless task. Yet under some assumptions – so-called causal completeness [153] – modelling domains is feasible. These assumptions state that the conditions provided are both necessary and sufficient for describing the action. (These concerns arise in causal modeling in machine learning as well [142], as we need to accurately identify all the parent variables that influence the variable of interest and describe them at the appropriate level of detail.)

If we do not make that assumption, the alternative approach would be to consider a wide range of typical cases, while also accounting for unusual and exceptional cases by incorporating the concept of abnormality. All of this requires notions of non-monotonicity.



Figure 5. Van Den Broeck on differentiability.

There is a general view that non-monotonicity is not needed and thereby a wasted effort, or already solved. Neal Parikh’s tweet, for example, has a view that non-monotonic reasoning is a wasted effort.

This seems to be a fairly superficial remark because there is no evidence that the problems identified in the non-monotonic reasoning community have been successfully addressed using any technique. It is true that many machine learning models, when trained on existing data, can identify typical patterns and detect abnormalities within that data [102,122]. However, there is no universal mechanism to address these concepts in a general way. Moreover, non-monotonic logic reasoning has given us notions like stable model semantics [72] which now powers recent approaches to neuro-symbolic learning [189]. Interestingly, non-monotonic semantics can also allow us to capture cycles in graphs [50], which ordinarily requires recursion using second-order logic [59]. This may be an important aspect as we utilize neural networks for reasoning about large graphs and the web more generally [136]. Thus, attempting to disregard this area of research seems premature.

3.7. Differentiability

Recent approaches to machine learning can be summarized by emphasizing the importance of differentiability as a key concept. However, it is widely held that logic cannot play a role in this. For example, Turing Award winner Yann LeCun quips [111]:

How can machine reason and plan in ways that are compatible with gradient-based learning?

Our best approaches to learning rely on estimating and using the gradient of a loss, which can only be performed with differentiable architectures and is difficult to reconcile with logic-based symbolic reasoning.

But as indicated by the sections above, this view is simply uniformed. Probabilities as well as real arithmetic can be mapped on to logical expressions and this means that both routes – a probabilistic one [69] and real-valued semantics [180] one – seem to naturally lead to differentiability. Let us elaborate further below.

There has been a historical understanding that logic and probability are compatible with each other [157,149,17]. These include topics such as 0-1 laws for studying the probability of satisfaction of first-order structures [61], the use of probability to compare the fit of logical hypothesis against observations [34], and perhaps most recently, the use of logic-based solvers by means of (weighted) model counting to compute conditional probabilities for Bayesian networks [38]. Consider the position in Guy Van Den Broeck’s tweet, for example.

At this point, there are plenty of approaches that explicitly use logic for the training of neural networks, especially in the context of regularization and differentiability. This started with the work of UCLA’s Semantic Loss [69] and KU Leuven’s DeepProbLog

[121], both of which adjust the loss function of the deep learning model based on a logical encoding of the constraints and program, respectively. This is an end-to-end approach in the sense that the predictions of the neural network are corrected using the logical solver and back-propagated to the network so that the trained network predicts outputs that are compatible with the constraints. There are also recent approaches that are based on real-valued variables, such as in [88] and in [180]. Providing arithmetic constraints to the training of deep learning networks and ensuring consistency with the provided domain knowledge is an important problem for areas like physics [167] and robotics [93].

However, it would be remiss not to point out that just because differentiability seems to be an important ingredient in the training of machine learning models, it does not mean that we expect every scientist in the area of logic to play game. There is still profound and rigorous work to be done on the integration of logical querying (e.g., computational effort needed to evaluate queries on a large knowledge basis [119]) and probability [13], for example. On the representation side, there are important issues to grapple with, such as languages to reason about logic and probability that permit the domain of quantification to be countably infinite (e.g., natural numbers) and uncountable (e.g., reals) sets [117]. Moreover, modal logics like temporal logics and dynamic logics become useful for deep learning-based endeavors as we navigate to more open-ended problems in dynamic domains [112]. For example, in [92], temporal logic formulas are used to train deep reinforcement learning agents. In [164,169], large language models are used to reason about dynamic epistemic properties [20], including the modelling of theory of mind [64]. And in [93], a temporally extended semantic loss function is considered.

An orthogonal direction of work that has recently been considered is the capturing of neural architectures, such as graph neural networks, using fragments of first-order logic [10]. For the purposes of our discussion, it suffices to say that simply focusing on differentiability or differentiable logics does not quite capture the range of questions that one can investigate in the AI landscape. Issues such as expressiveness, computational properties, and the development of hybrid architectures that combine the advantages of logical and uncertain reasoning continue to be valuable areas of research.

3.8. *What about “human-like” semantic definitions?*

The most well-studied semantics, or perhaps more accurately, the most widely-used semantics in computer science, remains classical [32]. That is, atoms are accorded values of either 0 or 1, and so formulas become Boolean functions. If modalities are introduced, such as time and actions [64], then we look at sequences of models: either a linear sequence or a tree-like sequence [153], for example.

But as mentioned above, there are also approaches where a degree of truth is accorded to formulas, either by allowing the atoms themselves to have non-binary values [190] or by according probabilities or other kinds of measures for complex formulas [56].

All of these notions are explored by establishing some kind of well-definedness, and logicians explore the implications of those conditions. For example, intuitionistic logic looks to weaken material implication [57]. Non-classical belief logics control the proof-depth of logical reasoners [118]. Fuzzy logic [190] was initially introduced with the idea that a truth definition needs to be provided to vague notions [66] such as being tall or making water warm.

Be that as it may, there is an informal argument often made that a mathematically rigorous definition of truth is too precise. Perhaps by training neural networks with real-world observations, they might exhibit more human-like reasoning capabilities that eschew a well-defined notion altogether? The evidence for this has not yet been established. Moreover, is such a feature desirable?

3.8.1. Flippant reasoning

Consider, for example, an informal inquiry by Oxford philosopher Luciano Floridi asks:⁷
what's the name of the daughter of Laura's mother?

ChatGPT (even GPT-4, the latest and best) – an influential large language model [187,27] – comes back with something to the effect of: “I’m sorry, I don’t have personal information about specific individuals.”

Analogously, the computer scientist Vincent Conitzer tried out the below puzzle:⁸

I put one apple in a bag. Then I put either one pear or two bananas in the bag. Then I remove one piece of fruit. Then I remove all the remaining pieces of fruit from the bag.

Based on what I told you, is it possible to know exactly how many pieces of fruit remain in the bag?

Here’s what we get:

Based on the information provided, it is not possible to know exactly how many pieces of fruit remain in the bag after removing one piece of fruit and then removing all remaining pieces. This is because there are two possible scenarios, but we don’t have enough information to determine which one took place:

- 1. One apple and one pear were in the bag. If you removed one piece of fruit, there would be one piece left. After removing all remaining pieces, there would be 0 pieces left.*
- 2. One apple and two bananas were in the bag. If you removed one piece of fruit, there would be two pieces left. After removing all remaining pieces, there would also be 0 pieces left.*

In both scenarios, after removing one piece of fruit and then removing all remaining pieces, there are no pieces of fruit left in the bag.

But apparently in a follow-up reply to that tweet, GPT-4 seems to get it correct.

What can we make of all of this? These models are clearly not reliable, or consistent [94]. Perhaps this what prompted Stephen Wolfram to consider a systematic integration of ChatGPT and Wolfram Alpha.⁹ By extension, a number of recent approaches seek to incorporate logical solvers as oracles [145] that can validate or disprove the predictions of neural architectures, such as large language models [191,140,139,127].

In sum, the “native” reasoning capabilities of purely neural models – that is, ones lacking a logic-based oracle – seem clearly limited.

3.8.2. Training vs deduction

An intriguing hypothesis put forward by eminent AI scientist Subbarao Kambhampati is this:¹⁰

⁷<https://twitter.com/Floridi/status/1635951391968567296?s=20>

⁸<https://twitter.com/conitzer/status/1636156048347111425?s=20>

⁹<https://writings.stephenwolfram.com/2023/03/chatgpt-gets-its-wolfram-superpowers/>

¹⁰<https://twitter.com/rao2z/status/1666294366720360449?s=20>

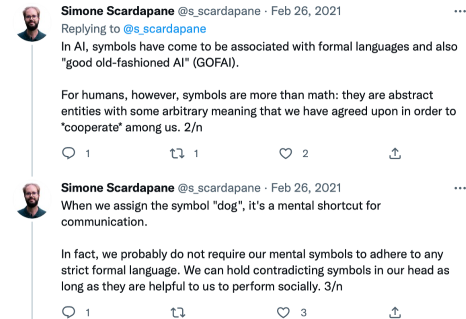


Figure 6. Scardapane on symbols and thought.

I think many of the claims about LLM’s reasoning capabilities miss the point that LLM’s are not just trained on “facts” but also, quite often, the deductive closure of those facts. Thus reasoning becomes (approximate) retrieval.

This hypothesis suggests that these models do not reason at all, but simply look for patterns of conclusions, which might limit, for example, the number of inference steps or the complexity of the reasoning process.

What about consoling ourselves with the idea that the training data might include all such deductions, in which case LLMs might be sufficient? Although work on such concerns continue, in a recent study [192], it is shown that LLMs likely pick up unnecessary statistical features of logical inputs, and their logical reasoning abilities may not be sound across different distributions on background theories, and thus, likely not complete.

3.8.3. *The intentional stance*

It is worth noting that, strictly speaking, we do not require that the semantics be given by humans, or that they be hand-written. Symbols can be obtained from low-level data (via symbol grounding), or from closely related languages [12], or from abstract descriptions [47] of concepts [109]. The use of symbols in AI also does not mean that symbolic logic experts assume humans manipulate symbols in their head. See [114] for philosophical discussions on this point, which can ultimately be tied to the “intentional stance” [51]. The intuition here is that any capability we attribute to an (artificial or human) agent could be understood in terms of intentions, beliefs, and other mental attitudes, which allow us to characterize what the agent is trying to do. It is a pragmatic perspective rather than a literal representation of the agent’s behavior model.

Consider the observation by Simone Scardapane in his tweet, for example. He suggests that the semantics of connectives and formulas may be built up from context, social environment and language use. While the search for a logic that accurately characterizes these kinds of observations with humans is still ongoing, it is worth noting that we do not necessarily need a logical knowledge basis to be consistent. For example, there is work on para-consistent logics [29].

Ultimately, we have a range of language choices to work with. We may disagree on the semantics, but having a few different systems that can be mathematically studied seems like a good start.

A follow-up question might be to the tune of: does it still make sense to bother with classical semantics? Just as it makes sense to study logic outside the context of differentiability, we would argue the study of classical semantics is also worthwhile in the AI context. Reasons include: (a) it is a well-defined mathematical model, (b) with the use of modalities and/or non-classical semantics, we can relate different systems, (c) we do not really know which semantics best approximates human reasoning, (d) we may not want mathematical truths that play fast and loose with inevitable conclusions just because we think humans might have some cognitive biases and exhibit inconsistent reasoning, and (e) the science of robust AI is still evolving. Logicians in AI should be allowed to investigate the properties of well-defined objects (including classical logic) with patience and rigour.

4. Logic and learning can be complementary

As already hinted above, symbolic logic can play an important role in training deep learning models but also in integrating reasoning as a post-hoc process or as a metalinguistic paradigm. That is, we can ensure that the distribution of the trained network respects domain constraints [88]. We can extract rules from trained models and reason about them outside the framework of the network [145]. Or we can use the outputs of the network as inputs to a computational paradigm such as probabilistic programming [121]. There is very interesting work on the semantics of programs that inherently support some notion of differentiation [1]. This is an object of intense theoretical study that can have consequences on the types of distributions that are expressible in programming languages [166]. So, this theory has far-reaching effects on what type of probabilistic models can be modelled effectively.

In the second half of the article, we make the following point: symbols and DL need not compete with each other, and can be complementary. Perhaps the most representative example of this is the burgeoning field of neuro-symbolic AI [70], which has come to encompass things like neural program induction [109], neural theorem improvers, and differentiable logics [191]. We consider some other categories below, as usual, with overlap.

4.1. Symbolic logic as meta-theory

An argument made previously [18] is that symbolic logic can be used to formalize notions currently out of the purview of standard machine learning. These include things like the semantics of involved probabilistic programming languages [166] and understanding the limits of differentiable logics [180], but it can also pertain to a range of more exotic topics.

For example, it is very common in AI applications these days to require frameworks for multi-agent reasoning [3]. In explainable AI [80], in particular, we might require that the robot holds beliefs about the human agent [99]. Modal logics study such phenomena. Thus, there has been a significant amount of recent work on incorporating agent modeling into learning frameworks, with multi-agent reinforcement learning being a prominent example [3]. Furthermore, incorporating agent modeling for explainable planning [2] and utilizing user-provided constraints as reward functions in reinforcement learning [92] are topics of study.

Moreover, complex AI systems are not going to be purely based on providing predictions. They will involve search, constraint reasoning, and planning [156]. This has necessitated new approaches for compositionality [166] and modularity [171]. In some recent work, for instance, it was noted that weighted model counting [74], which provides the foundation for mapping Bayesian inference to SAT solvers, can be upgraded to also reason about maximization and minimization of properties [101], leading to languages where a number of different AI sub-areas, such as search and optimization, can be unified [22].

An orthogonal but very interesting line of research in the recent years looks at the expressiveness of mainstream neural architectures using logical languages. Primarily, they look at fragments of first-order logic to capture (a simplified version of) neural architectures such as transformers [181] and graph neural networks [188]. These investigations have identified that graph neural networks capture fairly limited fragments of first-order logic [10], while attention mechanisms have been shown to be Turing-complete [144]. In the case of graph neural networks, the community is still exploring the implications of these results but it is believed that these architectures may fail in tasks involving queries that require more expressiveness than the fragment they correspond to. So, in this sense, using logical tools to understand neural architectures can have serious implications in terms of how these architectures are being used and in which circumstances they could be considered reliable.

4.2. High-level knowledge

At a number of recent AI events, Daniel Kahneman has been invited to discuss his famous distinction of the so-called *system 1* versus *system 2* type cognition in humans.¹¹ This is owing to the fact that AI scientists, for a very long time, have been deliberating on the appropriate way to abstract away low-level perception data with high-level concept knowledge, perhaps going back to Shakey [108]. Many “hybrid” formalisms for reasoning with perceptual data attempt to address the interplay between concepts and observations in a systematic way, e.g., [98,135].

Providing mechanisms as well as formal semantics for abstraction remains a topic of theoretical interest even today [14,89]. Roughly, the idea is given a representation R of the high-level model, to find another representation of R' involving low-level data and concepts, such that R and R' agree on atoms under a suitable mapping μ . That is, R entails an atom a iff R' entails $\mu(a)$. In a probabilistic setting, this might mean that we abstract a continuous distribution w.r.t. to evidence in terms of a discrete distribution [90]. There has also been some work on abstracting causal models [15].

In the specific case of deep learning systems, a key agenda point is how to define abstract concepts, whether extracted directly from data or defined externally, in order to coordinate and interoperate with these systems [109,21,33].

Consider the tweet by Gary Marcus, for example, further expanded in his recent book [123]. It is widely acknowledged that concepts such as time, abstraction, and causality will play a key role in designing an AI that has a world model that is rich enough to be interpreted in a way we would find reasonable [31]. Roughly, the idea is that given some system description, Σ , it is desirable to reason about the following:

¹¹<https://vimeo.com/390814190>

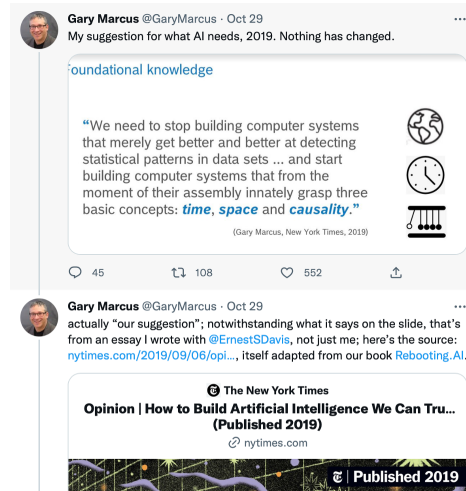
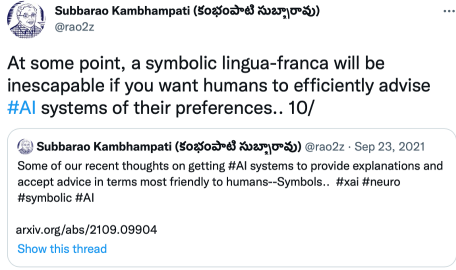


Figure 7. Marcus on reliable AI.

1. **Temporal abstractions:** Given two events, e and e' , we would like to know which happened earlier, and whether some trigger in e led to whatever happened in e' .
2. **Induction:** Suppose we have a set of events, p_1, \dots, p_n . We would like to find an idealized instance that generalizes these examples, \tilde{p} .
3. **Abstraction:** We would like find atomic descriptions, \hat{p} , that characterize the interaction between some of those instances (e.g., $p_1 \wedge p_2 \equiv \hat{p}$). The idea, then, is to use the abstract descriptor for increased comprehensibility [90,89].
4. **Causation:** Finally, given a causal chain from events p_1 to p_n in the sense that p_1 or its descendants causes p_n , we would like to understand what would happen if p_i was set to a certain value (*intervention*) or assumed a value not necessarily seen in the data (*counterfactual*).

Although there is some work on providing a causal semantics to deep learning systems [120], it is still in the early years and studied in a limited way. In contrast, we have very well-studied models of time [147] and causality with symbolic calculi [153,82,85]. It seems irresponsible to not utilise these frameworks simply because they are purely symbolic, and hence deemed “old-fashioned.”

As has been the case for many years now, symbols can be used as abstract identifiers for human-in-the-loop systems [99], and/or interactive machine learning especially when you have non-expert stakeholders engaging with predictors trained on high-dimensional data. See the position in Subbrao Kambhampati’s tweet, for example. In particular, there are very concrete examples from the neurosymbolic landscape that particularly highlight the benefits of using symbols. For example, the work on *reward machines* [92] looks to train deep learning-based reinforcement learning agents by means of high-level, temporally extended specifications, such as formulas expressed in linear temporal logic [37]. The propositions of the language are abstract descriptions of properties that can be understood by humans. There is also work on reasoning about neural concepts in a logical language. Although there have been prior works on hybrid formalisms that allow for machine learning constructs to be used in logic [98], recent neurosymbolic approaches



such as DeepProbLog [121] allow us to not only include neural concepts as objects in the logical program, but also to reason about this program as signals that could be fed back into the neural network training. This leads to a trained model that provides predictions and learns distributions that are consistent with the logical specification [87].

4.3. Symbolic logic can instantiate new methods of inference

One observation we emphasized earlier is that precisely because of the close relationship between logic and probability [34,177,16], it is possible to use logic-based solvers for doing probabilistic reasoning. This in turn, can mean that logic-based solvers are used in learning modules in probabilistic machine learning [178], or perhaps to reason about the output distributions of neural networks [69].

This is primarily instantiated via weighted model counting [74], which – as discussed above – is an extension of SAT solving to identify all possible satisfying assignments [6]. And as mentioned, there is also an extension of this strategy to deal with continuous properties via so-called weighted model integration [26]. One broader observation here is that because weighted model counting is defined in terms of weights on the possible models of a logical formula, it is possible to use different types of weights. This means a whole range of different computational tasks defined over the models of a logical formula can be approached using the same abstract specification of weighted model counting. This leads to the notion of *algebraic model counting* [101], where instead of sums over the models and products over the weights of literals, we can consider different kinds of corresponding operations such as *maximum* and *minimum* [6].

A notable development in this space is *knowledge compilation* [45]. This stems from the observation that given a probabilistic model, we may have to compute conditional queries repeatedly. Therefore, there have been efforts in representing a logical formula as a data structure that permits the computation of model counting [45], including in the presence of distinct conditional queries, effectively. This development can be coupled with the notion of algebraic model counting [101], but it has also served as a computational backbone for many emerging representations that unify logic and probability, such as relational Bayesian and Markov networks [179] – in addition to classical Bayesian networks [38], of course – and probabilistic logic programming languages like ProbLog [65].

Circuits provide a new way of doing inference with probabilistic models with the following properties: you pay a one-time cost for compiling the representation, such as a Bayesian network, into such a circuit, and then every query afterwards can be done in time polynomial in the size of the circuit. See the tweets by Kristian Kersting and



Figure 8. Kersting on circuits and hardware

Antonio Vergari, for example, which are representative of a broader program of learning such circuits directly [115]. The goal is to find an alternative to classical machine learning models with attractive computational properties for inference [182]. This is a new and exciting way of doing probabilistic reasoning and has even led to new approaches to inference in probabilistic programming [91].

4.4. Logical oracles

There is considerable work on verifying neural networks [162] for safety properties [35] as well as robustness [71], where we want to ensure that the prediction of neural networks does not change arbitrarily for small perturbations to the input. Along these lines, there is a new direction of work where logical reasoners serve as oracles to machine learning predictions to ensure that the predictions are consistent.

A representative example here is the contrasting of reasoning capabilities of large-scale learned models, such as large language models, against that of a symbolic oracle. Recent work on Wolfram Alpha [187] looks to integrate an arithmetic solver with the output of ChatGPT so that reasoning outputs are consistent and coherent with mathematical principles. Similarly, although there is some work on how the chain-of-thought prompting approach can lead to better reasoning outputs by large language models, the use of a logical oracle leads to provably correct outputs. The capabilities of ChatGPT, for example, have been directly studied in [67] and [94], and the use of a logical oracle to provide an externally sourced solution to reasoning problems with large language models is considered in [139]. In [169,164], such an approach has been shown to be applicable to involved problems involving the mental states of multiple agents, commonly referred to as the theory of mind [163,64].

In a rough sense, the idea here is not that different from the investigations on logic-based loss functions [69] because here too, predictions are expected to conform to logical constraints [87].

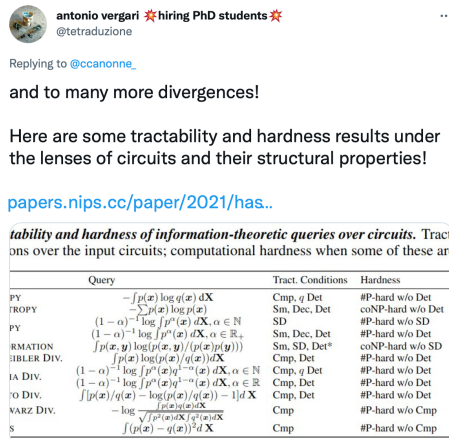


Figure 9. Vergari on information-theoretic properties of circuits in machine learning

4.5. Logic benefits from learning

In the article written so far, we have made the case for machine learning benefiting from logical tools and languages. However, on the other hand, looking back to the early days of logical thought, Aristotle argued for the importance of the process of induction [18]. We need mechanisms to learn the general from the particular, which involves generalizing from specific instances to create a generic statement that applies to all instances. That is, a quantified formula that entails all the atoms. In modern AI, this process is a key source of logical knowledge obtained from data [159,48], in addition to information provided by experts [46].

However, if our logical knowledge is to consist of a combination of expert-provided knowledge and knowledge drawn from examples, there are a number of concerns we need to address. For example, how can we ensure that a hypothesis that is consistent with the background knowledge is extracted from the observations [133]? What kind of properties should the resulting knowledge base have [128]? How do we deal with observations that might be incorrect or noisy [7]? How do we ensure that the formula we generalize from the observations captures not only the observations made so far but also the observations we have not yet seen and might encounter in the future [175,97]?

In recent years, a variety of approaches ranging from statistical relational learning [149] to probably-approximate correct (PAC) semantics [148] to neural program induction [109] and neural rule induction [60] have been explored. These approaches utilize state-of-the-art machine learning tools and theory to learn logical expressions. In some cases, noise in the observations is treated by assuming that the observations are drawn from an unknown distribution. In other cases, the generalization capabilities of neural networks are exploited to learn representations that are empirically robust to this noise.

It is now believed that machine learning will likely impact almost all of computer science because it provides a mechanism to construct models from data [161]. This means that we will continue considering combinations of model-based and data-driven domain knowledge in the future. All of this is even more reason to not entertain notions of dichotomy between logic and learning.

5. Concluding thoughts

In this article, we looked at a few of the misunderstandings that arise when considering the relevance and use of symbolic AI in modern AI systems. We have covered some of the ground that we feel frequently comes up. We hope the reader is convinced that not only do these dimensions have significant overlap – including ideas such as model counting appearing in and linking to multiple concerns – but it is also the case that recent advances are exploiting state-of-the-art machine learning, and in the process, improving on the state-of-the-art.

This speaks volumes in terms of why ignoring symbols and symbolic logic in general is a mistake and seems disingenuous from a scientific viewpoint, dishonest even. Whether there might be a future architecture that is very close in spirit to current neural models and makes logical tools redundant is yet to be seen. However, as we have argued, it is hard to imagine that, from a theoretical standpoint, logical analysis itself will become redundant, because many of the desired properties sought out are logical in nature. Despite reported advances in the reasoning capabilities of large language models, currently seen as the culmination of large-scale deep learning models, they still struggle with consistency and correctness in both logical and arithmetic problems.

5.1. Other dimensions

We have not discussed a few key issues that are emerging in the AI landscape. With the growing use of AI systems in financial and industrial applications, issues of trustworthiness and responsibility keep coming up [123].

For example, one area where symbolic logic is widely used in many stochastic systems [39] is the verification of safety properties [162], and/or testing for robustness [35]. The idea with safety properties is to ensure that certain regions in a geometric space are avoided because they might represent dangerous operational areas. In the case of robustness, we want to ensure that small perturbations to the input do not dramatically change the prediction from the neural network. It should not come as a surprise that ideas from logic-based computer science, including temporal logic [37] as well as satisfiability modulo theories [11], are the main tools to formalize and investigate these types of properties.

Another interesting avenue for examining trustworthy and responsible AI is understanding the ethical principles and norms under which AI systems should operate [54]. In this subarea, although mainstream models of concepts such as fairness do not necessarily use logic [183], further analysis of how systems could conform to ethical principles is often pursued through symbolic logic [52]. For example, notions such as act-deontology [105] or consequentialism can be formalized as properties that the system's execution should obey [138,186]. There has been work on using symbolic causal models to understand notions of blameworthiness, and the degree of responsibility [41]. Finally, there is considerable recent work on explainable planning [99], where a formal model is used to capture the user's intent and contrast it with the system's understanding of the world in which it operates [163]. For an overview on how knowledge representation can provide much needed frameworks for ethical and trustworthy AI, see [19].

5.2. *Neuro-symbolic AI*

As we discussed, one area where concerns about the use of logic seem to disappear is neuro-symbolic AI. Neuro-symbolic AI holds a lot of promise because it can offer interesting ways to combine symbolic logic and deep learning, and build on the success of both. And like the maxim: “the whole is greater than the sum of the parts,” such an integration may not simply be the communication of outputs in a divorced way, but could involve a deeper type of synthesis [86]. Some approaches have dealt with loss functions, while others have focused on post-hoc logical reasoning or extracting rules from networks. All of these approaches are interesting in their own right.

There is also a trade-off, at least as per our current understanding, between the complexity and level of detail of the logical knowledge and how effectively it can integrate with a learning system. For example, papers focusing on loss functions typically deal with smaller-sized formulas and constraints [88], while works exploring the integration of learning with knowledge graphs often consider ontologies with more than a hundred or even a thousand nodes [136]. Some may argue at this point whether these examples clearly indicate instances of neurosymbolic paradigms exceeding the capabilities of state-of-the-art machine learning. However, this is somewhat of a nebulous measure because state-of-the-art machine learning does encompass various neurosymbolic notions, even if they do not explicitly acknowledge it. Examples range from concept learning [109] to Wolfram Alpha-type integrations with large language models [187].

Of course, with such a diversity of solutions, it may be challenging to determine the correct approach. Perhaps there is no one-size-fits-all solution, and the combination of logic and deep learning can vary depending on the application. Regardless of the specific approach, it is clear that we need to understand the principles of logical languages and semantics to ensure that resulting mathematical objects are well-defined with desired properties. This appreciation is essential for both theoretical exploration and practical applications.

It should be noted that there is a case to be made for expressive representations. For example, some might come away feeling that the best way to approach the future of neuro-symbolic AI is to focus on very limited languages. But such a view may not be fruitful in the long term. For example, it is widely understood that first-order is useful for generalized assertions [113], and modal logics for time and multi-agent beliefs [64]. In general, the language is critical for capturing the domain correctly. In a statement remarkably similar in spirit, Judea Pearl writes [143]:

This is why you will find me emphasizing and reemphasizing notation, language, vocabulary and grammar. For example, I obsess over whether we can express a certain claim in a given language and whether one claim follows from others. My emphasis on language also comes from a deep conviction that language shapes our thoughts. You cannot answer a question that you cannot ask, and cannot ask a question that you have no words for.

And, as with Pearl and the knowledge representation community more generally, we will identify with “representation first, acquisition second.”

5.3. *No need to condescend*

To sum up, there is a lot to be gained by the relating the mathematical foundations of logic and deep learning. And the benefit is not purely for the logician, but also for the



Figure 10. Goldstaub’s commentary on Hinton

deep learning researcher who wants to think more broadly than prediction with big data. Scientists working on logic and language should be allowed to work on problems that seem scientifically relevant without necessarily linking to or competing with whatever the zeitgeist of machine learning is.

We should, of course, celebrate successes — its neither an accident nor misplaced opportunism that logic/programming language folks are interested in learning and are eager to understand the latest and best [79]. Moreover, what combination of logic and/or learning would be needed for general-purpose AI is not well-understood yet. We cannot point to the exact approach or balance of innateness vs tabula rasa we need for general AI, because we simply do not know. We can only loosely articulate requirements (e.g., correct, fair and safe by design), capabilities (e.g., ability to reason about causality, time and space models) and corresponding desiderata.

Experts can get excited about what works – the success of AlphaGo, as well as large language models, is kind of a success for AI, although of course it opens up questions about generality and correctness. However, there is no need to dismiss other approaches. Indeed, what we do not need are scientists – especially Turing award winners like Geoff Hinton – mocking other areas, such as the gasoline analogy. Likewise, we also do not need community members, such as Tabitha Goldstaub, with 16000+ followers on Twitter, sharing derision with conviction.¹²

Acknowledgements

This article builds and extends arguments previously presented in the author’s blog: medium.com/@vaishakbelle. The author is grateful to the Royal Society for funding his University Research Fellowship.

References

- [1] M. Abadi and G. D. Plotkin. A simple differentiable programming language. *Proceedings of the ACM on Programming Languages*, 4(POPL):1–28, 2019.

¹²<https://twitter.com/tabithagold/status/1070736319901519876?s=20>

- [2] S. V. Albrecht, C. Brewitt, J. Wilhelm, B. Gyevnar, F. Eiras, M. Dobre, and S. Ramamoorthy. Interpretable goal-based prediction and planning for autonomous driving. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1043–1049. IEEE, 2021.
- [3] S. V. Albrecht and P. Stone. Autonomous agents modelling other agents: A comprehensive survey and open problems. *Artificial Intelligence*, 258:66–95, 2018.
- [4] J. Andreas, M. Rohrbach, T. Darrell, and D. Klein. Neural module networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 39–48, 2016.
- [5] R. Azamfirei, S. R. Kudchadkar, and J. Fackler. Large language models and the perils of their hallucinations. *Critical Care*, 27(1):1–2, 2023.
- [6] F. Bacchus, S. Dalmao, and T. Pitassi. Solving #SAT and Bayesian inference with backtracking search. *J. Artif. Intell. Res. (JAIR)*, 34:391–442, 2009.
- [7] F. Bacchus, J. Y. Halpern, and H. J. Levesque. Reasoning about noisy sensors and effectors in the situation calculus. *Artificial Intelligence*, 111(1–2):171 – 208, 1999.
- [8] S. Badreddine, A. d. Garcez, L. Serafini, and M. Spranger. Logic tensor networks. *Artificial Intelligence*, 303:103649, 2022.
- [9] C. Baier and J.-P. Katoen. *Principles of Model Checking*. The MIT Press, 2008.
- [10] P. Barceló, E. V. Kostylev, M. Monet, J. Pérez, J. Reutter, and J.-P. Silva. The logical expressiveness of graph neural networks. In *8th International Conference on Learning Representations (ICLR 2020)*, 2020.
- [11] C. Barrett, R. Sebastiani, S. A. Seshia, and C. Tinelli. Satisfiability modulo theories. In *Handbook of Satisfiability*, chapter 26, pages 825–885. IOS Press, 2009.
- [12] S. Bartha, J. Cheney, and V. Belle. One down, 699 to go: or, synthesising compositional desugarings. *Proceedings of the ACM on Programming Languages*, 5(OOPSLA):1–29, 2021.
- [13] P. Beame, G. Van den Broeck, E. Gribkoff, and D. Suciú. Symmetric weighted first-order model counting. In *PODS*, pages 313–328. ACM, 2015.
- [14] S. Beckers and J. Y. Halpern. Abstracting causal models. In *Proceedings of the aaai conference on artificial intelligence*, volume 33, pages 2678–2685, 2019.
- [15] S. Beckers and J. Y. Halpern. Abstracting causal models. In *Proceedings of the aaai conference on artificial intelligence*, volume 33, pages 2678–2685, 2019.
- [16] V. Belle. Logic meets probability: Towards explainable ai systems for uncertain worlds. In *IJCAI*, 2017.
- [17] V. Belle. Symbolic logic meets machine learning: A brief survey in infinite domains. In *International Conference on Scalable Uncertainty Management*, pages 3–16. Springer, 2020.
- [18] V. Belle. Logic meets learning: From aristotle to neural networks. In *Neuro-Symbolic Artificial Intelligence: The State of the Art*, pages 78–102. IOS Press, 2021.
- [19] V. Belle. Knowledge representation and acquisition for ethical ai: challenges and opportunities. *Ethics and Information Technology*, 25(1):22, 2023.
- [20] V. Belle, T. Bolander, A. Herzig, and B. Nebel. Epistemic planning: Perspectives on the special issue, 2022.
- [21] V. Belle and A. Bueff. Deep inductive logic programming meets reinforcement learning. In *The 39th International Conference on Logic Programming*. Open Publishing Association, 2023.
- [22] V. Belle and L. De Raedt. Semiring programming: A semantic framework for generalized sum product problems. *International Journal of Approximate Reasoning*, 126:181–201, 2020.
- [23] V. Belle and G. Lakemeyer. Reasoning about probabilities in unbounded first-order dynamical domains. In *IJCAI*, 2017.
- [24] V. Belle, G. Lakemeyer, and H. J. Levesque. A first-order logic of probability and only knowing in unbounded domains. In *Proc. AAAI*, 2016.
- [25] V. Belle and H. J. Levesque. Reasoning about continuous uncertainty in the situation calculus. In *Proc. IJCAI*, 2013.
- [26] V. Belle, A. Passerini, and G. Van den Broeck. Probabilistic inference in hybrid domains by weighted model integration. In *IJCAI*, 2015.
- [27] A. Birhane, A. Kasirzadeh, D. Leslie, and S. Wachter. Science in the age of large language models. *Nature Reviews Physics*, pages 1–4, 2023.
- [28] S. Bistarelli, U. Montanari, and F. Rossi. Semiring-based constraint logic programming: syntax and semantics. *TOPLAS*, 23(1):1–29, 2001.
- [29] H. A. Blair and V. Subrahmanian. Paraconsistent logic programming. *Theoretical computer science*,

- 68(2):135–154, 1989.
- [30] R. Brachman and H. Levesque. *Knowledge representation and reasoning*. Morgan Kaufmann Pub, 2004.
 - [31] R. J. Brachman and H. J. Levesque. *Machines like us: toward AI with common sense*. MIT Press, 2022.
 - [32] A. R. Bradley and Z. Manna. *The calculus of computation: decision procedures with applications to verification*. Springer Science & Business Media, 2007.
 - [33] A. Bueff and V. Belle. Learning explanatory logical rules in non-linear domains: a neuro-symbolic approach. *Machine Learning*, pages 1–36, 2024.
 - [34] R. Carnap. *Logical foundations of probability*. Routledge and Kegan Paul London, 1951.
 - [35] M. Casadio, E. Komendantskaya, M. L. Daggitt, W. Kokke, G. Katz, G. Amir, and I. Refaeli. Neural network robustness as a verification property: a principled case study. In *International Conference on Computer Aided Verification*, pages 219–231. Springer, 2022.
 - [36] S. Chakraborty, D. J. Fremont, K. S. Meel, S. A. Seshia, and M. Y. Vardi. Distribution-aware sampling and weighted model counting for sat. *AAAI*, 2014.
 - [37] K. Chatterjee, M. Chmelik, R. Gupta, and A. Kanodia. Qualitative analysis of pomdps with temporal logic specifications for robotics applications. In *2015 IEEE International Conference on Robotics and Automation (ICRA)*, pages 325–330. IEEE, 2015.
 - [38] M. Chavira and A. Darwiche. On probabilistic inference by weighted model counting. *Artif. Intell.*, 172(6-7):772–799, 2008.
 - [39] T. Chen, M. Diciolla, M. Z. Kwiatkowska, and A. Mereacre. A simulink hybrid heart model for quantitative verification of cardiac pacemakers. In *Proceedings of the 16th international conference on Hybrid systems: computation and control, HSCC 2013, April 8-11, 2013, Philadelphia, PA, USA*, pages 131–136, 2013.
 - [40] D. Chistikov, R. Dimitrova, and R. Majumdar. Approximate counting in smt and value estimation for probabilistic programs. In *TACAS*, volume 9035, pages 320–334. 2015.
 - [41] H. Chockler and J. Y. Halpern. Responsibility and blame: A structural-model approach. *Journal of Artificial Intelligence Research*, 22:93–115, 2004.
 - [42] T. Crane. The language of thought: No syntax without semantics. 1990.
 - [43] A. Creswell, M. Shanahan, and I. Higgins. Selection-inference: Exploiting large language models for interpretable logical reasoning. *arXiv preprint arXiv:2205.09712*, 2022.
 - [44] A. Darwiche. Human-level intelligence or animal-like abilities? *Communications of the ACM*, 61(10):56–67, 2018.
 - [45] A. Darwiche, P. Marquis, D. Suciú, and S. Szeider. Recent trends in knowledge compilation (dagstuhl seminar 17381). In *Dagstuhl Reports*, volume 7. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2018.
 - [46] E. Davis. *Representations of commonsense knowledge*. Morgan Kaufmann, 2014.
 - [47] L. De Raedt. Logical settings for concept-learning. *Artificial Intelligence*, 95(1):187–201, 1997.
 - [48] L. De Raedt, A. Dries, I. Thon, G. Van den Broeck, and M. Verbeke. Inducing probabilistic relational rules from probabilistic examples. In *Proceedings of 24th international joint conference on artificial intelligence (IJCAI)*, volume 2015, pages 1835–1842. IJCAI-INT JOINT CONF ARTIF INTELL, 2015.
 - [49] L. De Raedt and A. Kimmig. Probabilistic (logic) programming concepts. *Machine Learning*, 100:5–47, 2015.
 - [50] M. Denecker, M. Bruynooghe, and V. Marek. Logic programming revisited: Logic programs as inductive definitions. *ACM Transactions on Computational Logic*, 2(4):623–654, 2001.
 - [51] D. C. Dennett. *The intentional stance*. MIT press, 1989.
 - [52] L. Dennis, M. Fisher, M. Slavkovik, and M. Webster. Formal verification of ethical choices in autonomous systems. *Robotics and Autonomous Systems*, 77:1–14, 2016.
 - [53] K. R. Dienes. String theory and the path to unification: A review of recent developments. *Physics Reports*, 287(6):447–525, 1997.
 - [54] V. Dignum. *Responsible artificial intelligence: how to develop and use AI in a responsible way*, volume 1. Springer, 2019.
 - [55] P. Domingos. *The master algorithm: How the quest for the ultimate learning machine will remake our world*. Basic Books, 2015.
 - [56] D. Dubois and H. Prade. Representation and combination of uncertainty with belief functions and possibility measures. *Computational Intelligence*, 4(3):244–264, 1988.

- [57] M. Dummett. The philosophical basis of intuitionistic logic. In *Studies in Logic and the Foundations of Mathematics*, volume 80, pages 5–40. Elsevier, 1975.
- [58] K. Ellis, A. Albright, A. Solar-Lezama, J. B. Tenenbaum, and T. J. O’Donnell. Synthesizing theories of human language with bayesian program induction. *Nature communications*, 13(1):5024, 2022.
- [59] H. Enderton. *A mathematical introduction to logic*. Academic press New York, 1972.
- [60] R. Evans and E. Grefenstette. Learning explanatory rules from noisy data. *Journal of Artificial Intelligence Research*, 61:1–64, 2018.
- [61] R. Fagin. Probabilities on finite models. *The Journal of Symbolic Logic*, 41(01):50–58, 1976.
- [62] R. Fagin and J. Y. Halpern. Reasoning about knowledge and probability. *J. ACM*, 41(2):340–367, 1994.
- [63] R. Fagin, J. Y. Halpern, and N. Megiddo. A logic for reasoning about probabilities. *Information and computation*, 87(1-2):78–128, 1990.
- [64] R. Fagin, J. Y. Halpern, Y. Moses, and M. Y. Vardi. *Reasoning About Knowledge*. MIT Press, 1995.
- [65] D. Fierens, G. V. den Broeck, I. Thon, B. Gutmann, and L. D. Raedt. Inference in probabilistic logic programs using weighted CNF’s. In *UAI*, pages 211–220, 2011.
- [66] K. Fine. Vagueness, truth and logic. 1997.
- [67] S. Frieder, L. Pinchetti, R.-R. Griffiths, T. Salvatori, T. Lukasiewicz, P. C. Petersen, A. Chevalier, and J. Berner. Mathematical capabilities of chatgpt. *arXiv preprint arXiv:2301.13867*, 2023.
- [68] H. Gaifman. Concerning measures in first order calculi. *Israel Journal of Mathematics*, 2(1):1–18, 1964.
- [69] K. Gajowniczek, Y. Liang, T. Friedman, T. Zabkowski, and G. Van den Broeck. Semantic and generalized entropy loss functions for semi-supervised deep learning. *Entropy*, 22(3):334, 2020.
- [70] A. S. d. Garcez, K. Broda, D. M. Gabbay, et al. *Neural-symbolic learning systems: foundations and applications*. Springer Science & Business Media, 2002.
- [71] T. Gehr, M. Mirman, D. Drachler-Cohen, P. Tsankov, S. Chaudhuri, and M. Vechev. Ai2: Safety and robustness certification of neural networks with abstract interpretation. In *2018 IEEE symposium on security and privacy (SP)*, pages 3–18. IEEE, 2018.
- [72] M. Gelfond and V. Lifschitz. The stable model semantics for logic programming. In *ICLP*, pages 1070–1080, 1988.
- [73] L. Getoor and B. Taskar. *Introduction to statistical relational learning (adaptive computation and machine learning)*. 2007.
- [74] C. P. Gomes, A. Sabharwal, and B. Selman. Model counting. In *Handbook of Satisfiability*. IOS Press, 2009.
- [75] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [76] N. D. Goodman, V. K. Mansingha, D. M. Roy, K. Bonawitz, and J. B. Tenenbaum. Church: a language for generative models. In *Proc. UAI*, pages 220–229, 2008.
- [77] K. Grace, J. Salvatier, A. Dafoe, B. Zhang, and O. Evans. When will ai exceed human performance? evidence from ai experts. *Journal of Artificial Intelligence Research*, 62:729–754, 2018.
- [78] T. P. Gros, H. Hermanns, J. Hoffmann, M. Klauck, and M. Steinmetz. Analyzing neural network behavior through deep statistical model checking. *International Journal on Software Tools for Technology Transfer*, 25(3):407–426, 2023.
- [79] S. Gulwani. Dimensions in program synthesis. In *PPDP*, pages 13–24. ACM, 2010.
- [80] D. Gunning. Explainable artificial intelligence (xai). Technical report, DARPA/I20, 2016.
- [81] J. Y. Halpern. *Reasoning about Uncertainty*. MIT Press, 2003.
- [82] J. Y. Halpern. *Actual causality*. MIT Press, 2016.
- [83] J. Y. Halpern and M. Y. Vardi. Model checking vs. theorem proving: A manifesto. In J. F. Allen, R. Fikes, and E. Sandewall, editors, *KR*, pages 325–334. Morgan Kaufmann, 1991.
- [84] C. S. Herrmann and M. Thielscher. Reasoning about continuous processes. In *AAAI/IAAI, Vol. 1*, pages 639–644, 1996.
- [85] C. Hitchcock. *Causality: Models, reasoning and inference*, 2001.
- [86] P. Hitzler. Neuro-symbolic artificial intelligence: The state of the art. 2022.
- [87] N. Hoernle, R. M. Karampatsis, V. Belle, and K. Gal. Multiplexnet: Towards fully satisfied logical constraints in neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 5700–5709, 2022.
- [88] N. Hoernle, R. M. Karampatsis, V. Belle, and K. Gal. Multiplexnet: Towards fully satisfied logical

- constraints in neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 5700–5709, 2022.
- [89] T. Hofmann and V. Belle. Abstracting noisy robot programs. *AAMAS*, 2023.
- [90] S. Holtzen, G. Broeck, and T. Millstein. Sound abstraction and decomposition of probabilistic programs. In *International Conference on Machine Learning*, pages 1999–2008. PMLR, 2018.
- [91] S. Holtzen, G. Van den Broeck, and T. Millstein. Scaling exact inference for discrete probabilistic programs. *Proceedings of the ACM on Programming Languages*, 4(OOPSLA):1–31, 2020.
- [92] R. T. Icarte, T. Q. Klassen, R. Valenzano, and S. A. McIlraith. Reward machines: Exploiting reward function structure in reinforcement learning. *Journal of Artificial Intelligence Research*, 73:173–208, 2022.
- [93] C. Innes and S. Ramamoorthy. Elaborating on learned demonstrations with temporal logic specifications. *arXiv preprint arXiv:2002.00784*, 2020.
- [94] M. Jang and T. Lukasiewicz. Consistency analysis of chatgpt. *arXiv preprint arXiv:2303.06273*, 2023.
- [95] E. T. Jaynes. How does the brain do plausible reasoning? In *Maximum-entropy and Bayesian methods in science and engineering: Foundations*, pages 1–24. Springer, 1988.
- [96] D. Jovanović and L. De Moura. Solving non-linear arithmetic. *ACM Communications in Computer Algebra*, 46(3/4):104–105, 2013.
- [97] B. Juba. Implicit learning of common sense for reasoning. In *Twenty-Third International Joint Conference on Artificial Intelligence*, 2013.
- [98] L. P. Kaelbling and T. Lozano-Pérez. Integrated task and motion planning in belief space. *I. J. Robotic Res.*, 32(9-10):1194–1227, 2013.
- [99] S. Kambhampati. Challenges of human-aware ai systems. *AI Magazine*, 41(3), 2020.
- [100] M. Katz. Łukasiewicz logic and the foundations of measurement. *Studia logica*, 40:209–225, 1981.
- [101] A. Kimmig, G. V. den Broeck, and L. D. Raedt. Algebraic model counting. *CoRR*, abs/1211.4475, 2012.
- [102] V. Kocijan, E. Davis, T. Lukasiewicz, G. Marcus, and L. Morgenstern. The defeat of the winograd schema challenge. *arXiv preprint arXiv:2201.02387*, 2022.
- [103] B. Konev, C. Lutz, A. Ozaki, and F. Wolter. Exact learning of lightweight description logic ontologies. *Journal of Machine Learning Research*, 18(201):1–63, 2018.
- [104] R. Kowalski and M. Sergot. A logic-based calculus of events. *New Generation Computing*, 4:67–95, 1986.
- [105] B. Krarup, F. Lindner, S. Krivic, and D. Long. Understanding a robot’s guiding ethical principles via automatically generated explanations. In *2022 IEEE 18th International Conference on Automation Science and Engineering (CASE)*, pages 627–632. IEEE, 2022.
- [106] S. Kripke. A completeness theorem in modal logic. *Journal of Symbolic Logic*, 24(1):1–14, 1959.
- [107] F. R. Kschischang, B. J. Frey, and H. Loeliger. Factor graphs and the sum-product algorithm. *IEEE Transactions on Information Theory*, 47(2):498–519, 2001.
- [108] B. Kuipers, E. A. Feigenbaum, P. E. Hart, and N. J. Nilsson. Shakey: from conception to history. *Ai Magazine*, 38(1):88–103, 2017.
- [109] B. M. Lake, R. Salakhutdinov, and J. B. Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- [110] G. Lakemeyer and H. J. Levesque. Cognitive robotics. In *Handbook of Knowledge Representation*, pages 869–886. Elsevier, 2007.
- [111] Y. LeCun. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27. *Open Review*, 62, 2022.
- [112] H. Levesque, R. Reiter, Y. Lespérance, F. Lin, and R. Scherl. Golog: A logic programming language for dynamic domains. *Journal of Logic Programming*, 31:59–84, 1997.
- [113] H. J. Levesque. *Thinking as Computation: A First Course*. MIT Press, 2012.
- [114] H. J. Levesque and G. Lakemeyer. *The logic of knowledge bases*. The MIT Press, 2001.
- [115] Y. Liang, J. Bekker, and G. Van den Broeck. Learning the structure of probabilistic sentential decision diagrams. In *Proceedings of the 33rd Conference on Uncertainty in Artificial Intelligence*, pages 134–145, 2017.
- [116] L. Libkin. *Elements of Finite Model Theory*. Springer, 2004.
- [117] D. Liu, Q. Feng, V. Belle, and G. Lakemeyer. Concerning measures in a first-order logic with actions and meta-beliefs. In *20th International Conference on Principles of Knowledge Representation and Reasoning*, 2023.

- [118] Y. Liu, G. Lakemeyer, and H. J. Levesque. A logic of limited belief for reasoning with disjunctive information. In *KR*, pages 587–597, 2004.
- [119] Y. Liu and H. J. Levesque. Tractable reasoning in first-order knowledge bases with disjunctive information. In *Proc. AAAI*, pages 639–644, 2005.
- [120] Y. Luo, J. Peng, and J. Ma. When causal inference meets deep learning. *Nature Machine Intelligence*, 2(8):426–427, 2020.
- [121] R. Manhaeve, S. Dumancic, A. Kimmig, T. Demeester, and L. De Raedt. Deepproblog: Neural probabilistic logic programming. *Advances in Neural Information Processing Systems*, 31, 2018.
- [122] G. Marcus. Am i human? *Scientific American*, 316(3):58–63, 2017.
- [123] G. Marcus and E. Davis. *Rebooting AI: Building artificial intelligence we can trust*. Vintage, 2019.
- [124] J. McCarthy. Programs with common sense. In *Proceedings of the Symposium on the Mechanization of Thought Processes*, National Physiology Lab, Teddington, England, 1958.
- [125] J. McCarthy. Applications of circumscription to formalizing common-sense knowledge. *Artificial Intelligence*, 28(1):89–116, 1986.
- [126] J. McCarthy and P. J. Hayes. Some philosophical problems from the standpoint of artificial intelligence. In *Machine Intelligence*, pages 463–502, 1969.
- [127] A. V. Miceli-Barone, F. Barez, I. Konstas, and S. B. Cohen. The larger they are, the harder they fail: Language models do not recognize identifier swaps in python. *arXiv preprint arXiv:2305.15507*, 2023.
- [128] L. Michael. Learning from partial observations. In *IJCAI*, pages 968–974. Citeseer, 2007.
- [129] B. Miich, B. Marthi, S. J. Russell, D. Sontag, D. L. Ong, and A. Kolobov. BLOG: Probabilistic models with unknown objects. In *Proc. IJCAI*, pages 1352–1359, 2005.
- [130] P. Minervini, M. Bosnjak, T. Rocktäschel, and S. Riedel. Towards neural theorem proving at scale. *arXiv preprint arXiv:1807.08204*, 2018.
- [131] P. P. Mitra. The circuit architecture of whole brains at the mesoscopic scale. *Neuron*, 83(6):1273–1283, 2014.
- [132] L. Morgenstern and S. A. McIlraith. John McCarthy’s legacy. *Artificial Intelligence*, 175(1):1 – 24, 2011.
- [133] S. Muggleton, L. De Raedt, D. Poole, I. Bratko, P. Flach, K. Inoue, and A. Srinivasan. Ilp turns 20. *Machine learning*, 86(1):3–23, 2012.
- [134] N. J. Nilsson. Probabilistic logic. *Artificial intelligence*, 28(1):71–87, 1986.
- [135] D. Nitti, V. Belle, T. De Laet, and L. De Raedt. Planning in hybrid relational mdps. *Machine Learning*, 106:1905–1932, 2017.
- [136] F. Niu, C. Zhang, C. Ré, and J. W. Shavlik. Deepdive: Web-scale knowledge-base construction using statistical learning and inference. *VLDS*, 12:25–28, 2012.
- [137] Z. Ognjanovic and M. Raškovic. Some first-order probability logics. *Theoretical Computer Science*, 247(1–2):191 – 212, 2000.
- [138] M. Pagnucco, D. Rajaratnam, R. Limarga, A. Nayak, and Y. Song. Epistemic reasoning for machine ethics with situation calculus. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 814–821, 2021.
- [139] L. Pan, A. Albalak, X. Wang, and W. Y. Wang. Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning. *arXiv preprint arXiv:2305.12295*, 2023.
- [140] D. Panas, S. Seth, and V. Belle. Can large language models put 2 and 2 together? probing for entailed arithmetical relationships. *arXiv preprint arXiv:2404.19432*, 2024.
- [141] J. Pearl. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Morgan Kaufmann, 1988.
- [142] J. Pearl. *Causality*. Cambridge university press, 2009.
- [143] J. Pearl and D. Mackenzie. *The book of why*.
- [144] J. Pérez, P. Barceló, and J. Marinkovic. Attention is turing-complete. *Journal of Machine Learning Research*, 22(75):1–35, 2021.
- [145] C. Persia and A. Ozaki. Extracting rules from neural networks with partial interpretations. *arXiv preprint arXiv:2204.00360*, 2022.
- [146] J. Prado, A. Chadha, and J. R. Booth. The brain network for deductive reasoning: a quantitative meta-analysis of 28 neuroimaging studies. *Journal of cognitive neuroscience*, 23(11):3483–3497, 2011.
- [147] A. Prior. *Past, present and future*. Oxford University Press, 1967.
- [148] A. Rader, I. G. Mocanu, V. Belle, and B. Juba. Learning implicitly with noisy data in linear arithmetic. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence*, 2021.

- [149] L. D. Raedt, K. Kersting, S. Natarajan, and D. Poole. Statistical relational artificial intelligence: Logic, probability, and computation. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 10(2):1–189, 2016.
- [150] L. D. Raedt, A. Kimmig, and H. Toivonen. Problog: A probabilistic prolog and its application in link discovery. In *Proc. IJCAI*, pages 2462–2467, 2007.
- [151] V. Raman, N. Piterman, and H. Kress-Gazit. Provably correct continuous control for high-level robot behaviors with actions of arbitrary execution durations. In *ICRA*, pages 4075–4081, Karlsruhe, Germany, 2013.
- [152] G. N. Reeke and G. M. Edelman. Real brains and artificial intelligence. *Daedalus*, pages 143–173, 1988.
- [153] R. Reiter. *Knowledge in action: logical foundations for specifying and implementing dynamical systems*. MIT Press, 2001.
- [154] M. Richardson and P. Domingos. Markov logic networks. *Machine learning*, 62(1):107–136, 2006.
- [155] J. Rintanen. Planning as satisfiability: Heuristics. *Artif. Intell.*, 193:45–86, 2012.
- [156] S. Russell and P. Norvig. *Artificial Intelligence: A Modern Approach*. Prentice-Hall, Englewood Cliffs, NJ, 2nd edition edition, 2003.
- [157] S. J. Russell. Unifying logic and probability. *Commun. ACM*, 58(7):88–97, 2015.
- [158] S. Sabour, N. Frosst, and G. E. Hinton. Dynamic routing between capsules. *Advances in neural information processing systems*, 30, 2017.
- [159] M. Sap, V. Schwartz, A. Bosselut, Y. Choi, and D. Roth. Commonsense reasoning for natural language processing. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts*, pages 27–33, 2020.
- [160] S. Sardina and Y. Lespérance. Golog speaks the bdi language. In *Programming Multi-Agent Systems*, volume 5919 of *LNCS*, pages 82–99. Springer Berlin Heidelberg, 2010.
- [161] R. B. Shapiro, R. Fiebrink, and P. Norvig. How machine learning impacts the undergraduate computing curriculum. *Communications of the ACM*, 61(11):27–29, 2018.
- [162] A. Shih, A. Darwiche, and A. Choi. Verifying binarized neural networks by angluin-style learning. In *Theory and Applications of Satisfiability Testing–SAT 2019: 22nd International Conference, SAT 2019, Lisbon, Portugal, July 9–12, 2019, Proceedings 22*, pages 354–370. Springer, 2019.
- [163] M. Shvo, T. Q. Klassen, and S. A. McIlraith. Towards the role of theory of mind in explanation. In *Explainable, Transparent Autonomous Agents and Multi-Agent Systems: Second International Workshop, EXTRAAMAS 2020, Auckland, New Zealand, May 9–13, 2020, Revised Selected Papers 2*, pages 75–93. Springer, 2020.
- [164] D. Sileo and A. Lernould. Mindgames: Targeting theory of mind in large language models with dynamic epistemic modal logic. *arXiv preprint arXiv:2305.03353*, 2023.
- [165] P. Smolensky. Connectionist ai, symbolic ai, and the brain. *Artificial Intelligence Review*, 1(2):95–109, 1987.
- [166] S. Staton, H. Yang, F. Wood, C. Heunen, and O. Kammar. Semantics for probabilistic programming: higher-order functions, continuous distributions, and soft constraints. In *Proceedings of the 31st Annual ACM/IEEE Symposium on Logic in Computer Science*, pages 525–534, 2016.
- [167] R. Stewart and S. Ermon. Label-free supervision of neural networks with physics and domain knowledge. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- [168] L. W. Swanson. *Brain architecture: understanding the basic plan*. Oxford University Press, USA, 2012.
- [169] W. Tang and V. Belle. Tom-lm: Delegating theory of mind reasoning to external symbolic executors in large language models. *arXiv preprint arXiv:2404.15515*, 2024.
- [170] S. Tellex, T. Kollar, S. Dickerson, M. R. Walter, A. G. Banerjee, S. Teller, and N. Roy. Approaching the symbol grounding problem with probabilistic graphical models. *AI magazine*, 32(4):64–76, 2011.
- [171] E. Ternovska and D. G. Mitchell. Declarative programming of search problems with built-in arithmetic. In *Proc. IJCAI*, pages 942–947, 2009.
- [172] S. Thrun, W. Burgard, and D. Fox. *Probabilistic Robotics*. MIT Press, 2005.
- [173] J. Tooby, L. Cosmides, and H. C. Barrett. Resolving the debate on innate ideas. *The innate mind: Structure and content*, pages 305–337, 2005.
- [174] A. M. Turing. Computing machinery and intelligence. *Mind*, 59(236):433–460, 1950.
- [175] L. G. Valiant. Robust logics. In *Proceedings of the thirty-first annual ACM symposium on Theory of Computing*, pages 642–651, 1999.

- [176] K. Valmeekam, A. Olmo, S. Sreedharan, and S. Kambhampati. Large language models still can't plan (a benchmark for llms on planning and reasoning about change). *arXiv preprint arXiv:2206.10498*, 2022.
- [177] J. van Benthem. Against all odds: when logic meets probability. *ModelEd, TestEd, TrustEd: Essays Dedicated to Ed Brinksma on the Occasion of His 60th Birthday*, pages 239–253, 2017.
- [178] G. Van den Broeck. *Lifted Inference and Learning in Statistical Relational Models*. PhD thesis, KU Leuven, 2013.
- [179] G. Van den Broeck, W. Meert, and J. Davis. Lifted generative parameter learning. In *Statistical Relational Artificial Intelligence, AAAI Workshop*, 2013.
- [180] E. van Krieken, E. Acar, and F. van Harmelen. Analyzing differentiable fuzzy logic operators. *Artificial Intelligence*, 302:103602, 2022.
- [181] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [182] A. Vergari, Y. Choi, A. Liu, S. Teso, and G. Van den Broeck. A compositional atlas of tractable circuit operations for probabilistic inference. *Advances in Neural Information Processing Systems*, 34:13189–13201, 2021.
- [183] S. Verma and J. Rubin. Fairness definitions explained. In *2018 IEEE/ACM International Workshop on Software Fairness (Fairware)*, pages 1–7. IEEE, 2018.
- [184] K. Wang, E. Tsamoura, and D. Roth. On learning latent models with multi-instance weak supervision. *Advances in Neural Information Processing Systems*, 36, 2024.
- [185] P.-W. Wang, P. Donti, B. Wilder, and Z. Kolter. Satnet: Bridging deep learning and logical reasoning using a differentiable satisfiability solver. In *International Conference on Machine Learning*, pages 6545–6554. PMLR, 2019.
- [186] A. Winfield, C. Blum, and W. Liu. Towards an ethical robot: internal models, consequences and ethical action selection. In *Proceedings of the 14th Conference Towards Autonomous Robotic Systems*, pages 85–96, 2014.
- [187] S. Wolfram. Wolfram— alpha as the way to bring computational knowledge superpowers to chatgpt. *Stephen Wolfram Writings RSS, Stephen Wolfram, LLC*, 9, 2023.
- [188] K. Xu, W. Hu, J. Leskovec, and S. Jegelka. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826*, 2018.
- [189] Z. Yang, A. Ishay, and J. Lee. Neurasp: Embracing neural networks into answer set programming. In *29th International Joint Conference on Artificial Intelligence (IJCAI 2020)*, 2020.
- [190] L. A. Zadeh. Fuzzy sets. *Information and control*, 8(3):338–353, 1965.
- [191] H. Zhang, J. Huang, Z. Li, M. Naik, and E. Xing. Improved logical reasoning of language models via differentiable symbolic programming. *arXiv preprint arXiv:2305.03742*, 2023.
- [192] H. Zhang, L. H. Li, T. Meng, K.-W. Chang, and G. V. d. Broeck. On the paradox of learning to reason from data. *arXiv preprint arXiv:2205.11502*, 2022.