# From Statistical Relational to Neurosymbolic Artificial Intelligence: a Survey.

Giuseppe Marra[a], Sebastijan Dumančić[c], Robin Manhaeve[a], Luc De Raedt[a,b]

[a]*KU Leuven, Department of Computer Science and Leuven.AI*
[b]*Örebro University, Center for Applied Autonomous Sensor Systems*
[c]*Delft University of Technology, Department of Software Technology*

## Abstract

This survey explores the integration of learning and reasoning in two different fields of artificial intelligence: neurosymbolic and statistical relational artificial intelligence. Neurosymbolic artificial intelligence (NeSy) studies the integration of symbolic reasoning and neural networks, while statistical relational artificial intelligence (StarAI) focuses on integrating logic with probabilistic graphical models. This survey identifies seven shared dimensions between these two subfields of AI. These dimensions can be used to characterize different NeSy and StarAI systems. They are concerned with (1) the approach to logical inference, whether model or proof-based; (2) the syntax of the used logical theories; (3) the logical semantics of the systems and their extensions to facilitate learning; (4) the scope of learning, encompassing either parameter or structure learning; (5) the presence of symbolic and subsymbolic representations; (6) the degree to which systems capture the original logic, probabilistic, and neural paradigms; and (7) the classes of learning tasks the systems are applied to. By positioning various NeSy and StarAI systems along these dimensions and pointing out similarities and differences between them, this survey contributes fundamental concepts for understanding the integration of learning and reasoning.

## 1. Introduction

The integration of learning and reasoning is a key challenge in artificial intelligence and machine learning today. Various communities are addressing it, especially the field of neurosymbolic artificial intelligence (NeSy) [12, 28]. NeSy's goal is to integrate symbolic reasoning with neural networks. NeSy already has a long tradition, and it has recently attracted a lot of attention. Indeed, the topic has been addressed by prominent researchers such as Y. Bengio and H. Kautz in their keynotes at AAAI 2020, by Y. Bengia and G. Marcus in the AI Debate

---

[10] and Hochreiter has recently stated [56] that NeSy is "the most promising approach to a broad AI".

Another domain with a rich tradition in integrating learning and reasoning is that of statistical relational learning and artificial intelligence (StarAI) [44, 96]. StarAI focuses on integrating logical and probabilistic reasoning.

Historically, these two endeavours have adopted different learning paradigms, probabilistic versus neural, for integrating logic into machine learning. This in turn has resulted in two different subcommunities. StarAI has focused on probabilistic logics, their semantics and making inference more tractable, while learning is usually based on parameter learning techniques from probabilistic graphical models. On the other hand, NeSy extends neural networks with symbolic knowledge, focusing on scalable approximate models, paying less attention to semantical issues. In particular, NeSy techniques can often be characterized by a clear parameterization in terms of neural networks, i.e. layered structures of latent representations, and by resorting to the gradient-based backpropagation paradigm for learning. Despite a different focus and approach, the two domains want to achieve the same goal, that is, to integrate learning and reasoning. It is therefore surprising that there has been relatively little interaction between the two domains, but see [27, 11].

This discrepancy is the key motivation behind this survey: it aims at pointing out the similarities between these two endeavours and, in this way, it wants to stimulate cross-fertilization. We start from the literature on StarAI, following the key concepts, definitions and techniques outlined in several textbooks and tutorials such as [104, 96], because it turns out that the same issues and techniques that arise in StarAI apply to NeSy as well.

The key contributions of this paper are:

1. *We identify seven dimensions that these fields have in common and that can be used to categorize both StarAI and NeSy approaches.* These seven dimensions are concerned with (1) model vs proof-based inference, (2) logic syntax, (3) semantics, (4) learning parameters or structure, (5) representing entities as symbols or subsymbols, (6) integrating logic with probabilistic and/or neural concepts, and (7) learning tasks.

2. We provide evidence for our claim by positioning a wide range of StarAI and NeSy systems along these dimensions and pointing out analogies between them. This provides not only new insights into the relationships between StarAI and NeSy, but it also allows one to carry over and adapt techniques from one field to another. These insights provide opportunities for cross-fertilization between StarAI and NeSy, by focusing on those dimensions that have not been fully exploited yet.

3. We gently introduce key logical concepts and techniques inherited from StarAI. In this way, the paper also provides a gentle introduction to symbolic AI and StarAI techniques for the interested "connectionist" practitioner.

4. We illustrate each dimension using existing methods, and in this way, also present an intuitive and concrete overview of the research field.

Unlike some other perspectives on neurosymbolic computation [12, 28, 16], the present survey limits itself to a logical perspective and to developments in neurosymbolic computation that are consistent with this perspective. Therefore, we usually refer to symbols and symbolic algorithms as synonyms for logical representations and logical reasoning. Furthermore, the survey focuses on representative and prototypical systems rather than aiming at completeness (which would not be possible given the fast developments in the field). Several other surveys about neurosymbolic AI have been proposed. An early overview of neurosymbolic computation is that of [4]. Unlike the present survey it focuses very much on a logical and a reasoning perspective. Today, the focus has shifted very much to learning. More recently, [67] analysed the intersection between NeSy and graph neural networks (GNN). [123] described neurosymbolic systems in terms of the composition of blocks described by few patterns, concerning processes and exchanged data. In contrast, this survey is more focused on the underlying principles that govern such a composition. Finally, [26] exploits a neural network viewpoint by investigating in which components (i.e. input, loss or structure) symbolic knowledge is injected.

*Structure of the paper.* The next seven sections each describe one dimension by first introducing the main underlying concepts, either based on logic, probability or machine learning, and then showing how they are incorporated in StarAI and NeSy systems. Section 2 presents how to use logic for inference by distinguishing between proof-based and model-based systems, while Section 3 introduces logic at the syntac level, in particular, propositional, relational and first-order logic. Section 4 then introduces the semantics of logic and shows how to extend it to a continuous semantics, using fuzzy and probabilistic logics. Section 5 discusses the dimension of learning, distinguishing parameter learning from structure learning. Section 6 focuses on the representational level and to what extent neurosymbolic models use symbolic and/or subsymbolic features. Section 7 positions neurosymbolic approaches along the spectrum of three main paradigms, i.e. logic, probability and neural networks. Section 8 describes general classes of learning tasks to which neurosymbolic systems are usually applied. Finally, in Section 10, we conclude by introducing open challenges in the neurosymbolic landscape.

We summarize various neurosymbolic approaches along these dimensions in Table 1.

| Frameworks | Inference (P)roof (M)odel | Syntax (P)ropositional (R)elational (FOL) | Semantics (M)inimal (S)table (C)lassical (F)uzzy (P)robability | Learning (P)arameters (S)tructure | Representations (S)ymbolic (Sub)symbolic | Paradigms Logic (L/l) Probability (P/p) Neural(N/n) | Tasks (D)istant (S)upervision (S)emi (S)upervised (KGC)ompletion (G)enerative (K)nowledge (I)nduction |
|---|---|---|---|---|---|---|---|
| αILP [111] | P+M | FOL | S + P | P + S | S | Ln | KI |
| ∂ILP [39] | P | R | M + F | P + S | S | Ln | DS + KI |
| DeepProbLog [72] | P+M | FOL | M + P | P+S | S+Sub | LpN | DS + KI |
| DeepStochLog[132] | P | FOL | M + P | P | S | LpN | DS + SS |
| DiffLog[112] | P | R | M + F | P+S | S | Ln | KI |
| DL2[40] | M | P | C + F | P | S+Sub | IN | DS + SS |
| DLM[77] | M | FOL | C + F + P | P | S | IPN | SS + KGC |
| LRNN[116] | P | R | M + F | P + S | S + Sub | LN | KGC + KI |
| LTN[5] | M | FOL | C + F | P | S + Sub | IN | DS + SS |
| NeuralLP[137] | P | R | M + F | P | S | Ln | KGC + KI |
| NeurASP[138] | P+M | FOL | S + P | P | S | LpN | DS |
| NLM[35] | P | R | M + F | P + S | S | Ln | KGC + KI |
| NLog[121] | P | R | M + P | P | S | LpN | DS |
| NLProlog[131] | P | R | M + P | P + S | S + Sub | LpN | KGC + KI |
| NMLN[78] | M | FOL | C + P | P + S | S + Sub | IPN | KGC + G |
| NTP[102] | P | R | M + F | P + S | S + Sub | Ln | KGC + KI |
| RNM[76] | M | FOL | C + P | P | S | IPN | SS |
| SBR[33] | M | FOL | C + F | P | S+Sub | IN | DS + SS |
| Scallop[58] | P | FOL | M + P | P | S | LpN | DS |
| SL[133] | M | P | C + P | P | S | LpN | SS |
| Slash[113] | P+M | FOL | S + P | P | S | LpN | DS +SS |
| TensorLog[18] | P | R | M + P | P | S | LpN | DS + KGC |

Table 1: Logic-based NeSy frameworks according to the 6 dimensions outlined in the paper.

## 2. Proof- vs Model-theoretic View of Logic

In this paper, we focus on clausal logic as it is a standard form to which any first order logical formula can be converted. In clausal logic, theories are represented in terms of clauses. More formally, a *clause* is an expression of the form $h_1 \vee ... \vee h_k \leftarrow b_1 \wedge ... \wedge b_n$. The $h_k$ are *head literals* or conclusions, while the $b_i$ are *body literals* or conditions. Clauses with no conditions ($n = 0$) and one conclusion ($k = 1$) are *facts*. Clauses with only one conclusion ($k = 1$) are *definite clauses*.

The question we want to answer in this section is how to use such clausal theories to reason? And, how to infer new facts from the known clauses? Along this first dimension, we will investigate the two fundamental ways to view logical inference and determine the implications for StarAI and NeSy systems. In one view, we want to find *proofs* for a certain query, which leads to the proof-theoretic approach to logic. In the other view, we want to find *models* (that is, truth assignments to the logical atoms) that satisfy a given theory. This leads to the model-theoretic approach to logic.

*Proof-theoretic logic.* The proof-theoretic approach finds proofs for a query in a logic theory. While this approach to inference is applicable to any logic theory, we focus on logic programs in this paper. Syntactically, a logic program is a *definite* clause theory, which is a theory where all the clauses are definite (i.e. only one conclusion). In logic programs, definite clauses are interpreted as if-then **rules** ($h$ is true if $b_1, ..., b_n$ are true).

A proof for a query $q$ is a sequence of logical inference steps that demonstrates the truth of a query based on the given program. A compact way of representing the set of all proofs in a logic program uses an AND/OR tree, which consists of AND and OR nodes and edges amongst them. Each node represents a goal. An AND node branches into one or more outgoing edges, each representing one of the sub-goals that need to be simultaneously satisfied for the goal in the AND node to be true. An OR node represents choices or alternatives between multiple clauses that can be used to prove a particular sub-goal. The OR node branches into multiple outgoing edges, each representing one of these possible choices. Leaf nodes in the AND/OR tree represent true facts. Typically, forward or backward chaining inference is used to search for proofs for queries. We illustrate this in Example 1.

---

EXAMPLE 1: LOGIC PROGRAMS AND PROOFS

Consider the following logic program (adapted from [91]):

```
burglary.
hears_alarm_mary.

earthquake.
hears_alarm_john.
```
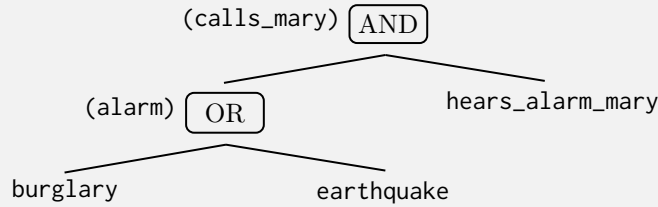
---

```
                  alarm <- earthquake.
                  alarm <- burglary.
                  calls_mary <- alarm,hears_alarm_mary.
                  calls_john <- alarm,hears_alarm_john.
```

The rules for *alarm* state that there will be an *alarm* if there is a *burglary* or an *earthquake*.
The set of proofs for the query *calls_ mary* can be represented compactly as an AND/OR tree.



*Model-theoretic logic.* On the other hand, the model theoretic perspective on logic is to find a model or truth assignment to the logical atoms that satisfy a given logic theory. An *interpretation*, or possible world, is a truth-assignment to the propositions (or ground atoms) of the language, and can be uniquely identified with the set of propositions it assigns $True$ (thus considering all the other $False$). An interpretation is a *model* of a clause $h_1 \vee ... \vee h_k \leftarrow b_1 \wedge ... \wedge b_n$ if at least one of the $h_i$ is in the interpretation when all the $b_1 \wedge ... \wedge b_n$ are in the interpretation as well. An interpretation $I$ is a model of a theory $T$, and we write $I \models T$, if it is a model of all clauses in the theory. We say that the theory is *satisfiable* if it has a model. The satisfiability problem, that is, deciding whether a theory has a model, is one of the most fundamental ones in computer science (cf. the SAT problem for propositional logic).

EXAMPLE 2: MODEL-THEORETIC

Consider the following theory (adopted from [100]):

$$smokes\_mary \leftarrow smokes\_john, influences\_john\_mary.$$
$$smokes\_john \leftarrow smokes\_mary, influences\_mary\_john.$$

$$smokes\_mary \leftarrow stress\_mary.$$
$$smokes\_john \leftarrow stress\_john.$$

A model of the previous theory is the set:

$$M = \{\mathtt{stress\_john}, \mathtt{smokes\_john}\}$$

By considering all the elements of this set *True* and all others *False*, the four clauses are satisfied.

In the model-theoretic perspective, one uses the logic theory as a set of *constraints* on the propositions, that is, the propositions are related to one another, without imposing a directed inference relationship between them as in forward or backward chaining. More details on these connections can be found in [96, 41].

### 2.1. Implications for StarAI

Statistical Relational AI's focus is on unifying logical and probabilistic graphical models (PGMs). A PGM [66] is a graphical model that compactly represents a (joint) probability distribution $P(X_1, ..., X_n)$ over $n$ discrete or continuous random variables $X_1, ..., X_n$. The key idea is that the joint factorizes over some factors $f_i$ specified over subsets $X^i$ of the variables $\{X_1, ..., X_n\}$.

$$P(X_1, ..., X_n) = \frac{1}{Z} f_1(X^1) \times ... \times f_k(X^k)$$

The random variables correspond to the nodes in the graphical structure, and the factorization is determined by the edges in the graph.

There are two classes of graphical models: *directed* ones, or Bayesian networks, and *undirected* ones, or Markov Networks. In Bayesian networks, the underlying graph structure is a directed acyclic graph, and the factors $f^i(X_i|parents(X_i))$ correspond to the conditional probabilities $P(X_i|parents(X_i))$, where $parents(X_i)$ denotes the set of random variables that are a parent of $X_i$ in the graph. In Markov networks, the graph is undirected and the factors $f^i(X^i)$ correspond to the set of nodes $X^i$ that form (maximal) cliques in the graph. Furthermore, the factors are non-negative and $Z$ is a normalisation constant.

The distinction between directed and undirected graphical models is parallel to the proof- vs model-theoretic view of logic. This parallel is at the very core of StarAI. In fact, by viewing each variable $X_i$ (or proposition) *at the same time* as a random and as a logical variable [105], clausal theories can be extended to define probabilistic models. Clauses can then be translated into binary valued factors by labeling them with weights (or probabilities), thus parameterizing the corresponding factors.

In the remainder of this section, we will show how StarAI has used this parallel to define two types of systems [96].

The first type of StarAI system generalizes directed models and resembles Bayesian networks. It includes well-known representations such as plate notation [66], probabilistic relational models (PRMs) [43], probabilistic logic programs (PLPs) [97], and Bayesian logic programs (BLPs) [61]. Today the most typical and popular representatives of this category are the probabilistic (logic) programs.

Probabilistic logic programs were introduced by Poole [92] and the first learning algorithm is due to Sato [105]. Probabilistic logic programs are essentially definite clause programs where every fact is annotated with the probability that it is *True*. This then results in a possible world semantics. The reason why probabilistic logic programs are viewed as directed models is clear when
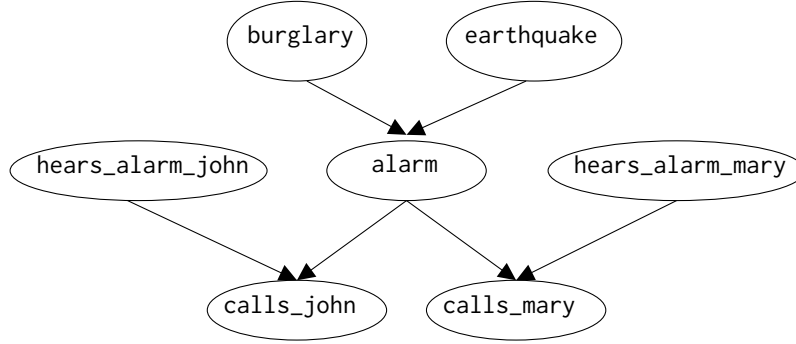
Figure 1: The Bayesian network corresponding to the ProbLog program in Example 3

looking at the derivations for a query, cf. Example 1. At the top of the AND-OR tree, there is the query that one wants to prove and the structure of the tree is that of a directed graph (even though it need not be acyclic). One can straightforwardly map directed graphical models, that is, Bayesian networks, onto such probabilistic logic programs by associating one definite clause to every entry in the conditional probability tables, yielding factors of the form $P(X|Y_1, ..., Y_n)$. Assuming boolean random variables, each entry $x, y_1, ..., y_n$ with parameter value $v$ can be represented using the definite clause $X(x) \leftarrow Y_1(y_1) \land ... \land Y_n(y_n) \land p_{x,y_1,...,y_n}$ and probabilistic fact $v :: p_{x,y_1,...,y_n}$. A probabilistic version of Example 1 is shown in Example 3 using the syntax of ProbLog [31].

---

EXAMPLE 3: PROBLOG

We show a probabilistic extension of the alarm program using ProbLog.

```
0.1::burglary.
0.3::hears_alarm_mary.
0.05::earthquake.
0.6::hears_alarm_john.

alarm <- earthquake.
alarm <- burglary.
calls_mary <- alarm,hears_alarm_mary.
calls_john <- alarm,hears_alarm_john.
```

This program can be mapped to the Bayesian network in Figure 1
This probabilistic logic program defines a distribution $p$ over possible worlds $\omega$. Let $P$ be a Problog program and $F = \{p_1 :: c_1, \cdots, p_n :: c_n\}$ be the set of ground probabilistic facts $c_i$ of the program and $p_i$ their corresponding probabilities. ProbLog defines a probability distribution over $\omega$ in the following way:
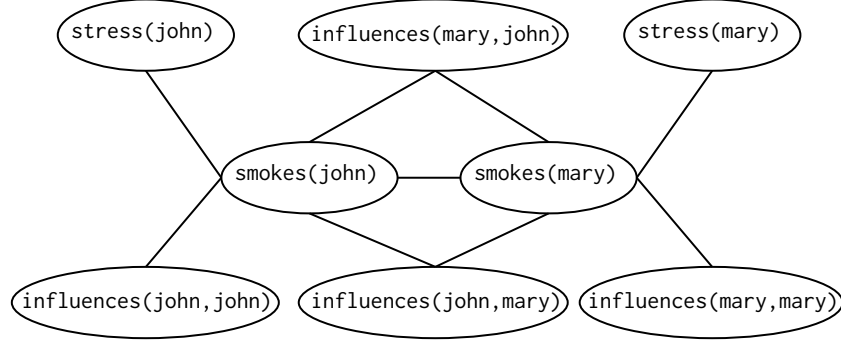
Figure 2: The Markov Field corresponding to the Markov logic network in Example 4

$$p(\omega) = \begin{cases} 0, & \text{if } \omega \not\models P \\ \displaystyle\prod_{c_i \in \omega : c_i = T} p_i \cdot \prod_{c_j \in \omega : c_j = F} (1 - p_j), & \text{if } \omega \models P \end{cases}$$

The second type of StarAI system generalizes undirected graphical models such as Markov networks or random fields. The prototypical example is Markov Logic Networks (MLNs) [100], and also Probabilistic Soft Logic (PSL) [3] follows this idea.

Undirected StarAI models consist of a set of weighted clauses $w : h_1 \vee ... \vee h_k \leftarrow b_1 \wedge ... \wedge b_m$ that become soft constraints. The higher the weight of a ground clause, the less likely possible worlds that violate these constraints are. In the limit, when the weight is $+\infty$ the constraint must be satisfied and becomes a purely logical constraint, a hard constraint. The weighted clauses specify a more general relationship between the conclusion and the condition than the definite clauses of directed models. While clauses of undirected models can still be used in (resolution) theorem provers, they are commonly viewed as constraints that relate these two sets of atoms.

Such undirected StarAI models can be mapped to an undirected probabilistic graphical model in which there is a one-to-one correspondence between grounded weighted clauses and factors, as we show in Example 4.

EXAMPLE 4: MARKOV LOGIC NETWORKS

We show a probabilistic extension (adapted from [100]) of the theory in Example 2 using the formalism of Markov Logic Networks. We use a First Order language with domain $D = \{john, mary\}$ and weighted clauses $\alpha_1$ and

9

$\alpha_2$, i.e.:

$$\alpha_1 : \quad 2.0 :: \texttt{smokes(Y)} \leftarrow \texttt{smokes(X)}, \texttt{influences(X,Y)}$$
$$\alpha_2 : \quad 0.5 :: \texttt{smokes(X)} \leftarrow \texttt{stress(X)}$$

In Figure 2, we show the corresponding Markov field.

A Markov Logic Network defines a probability distribution over possible worlds as follows. Let $A = [\alpha_1, \cdots, \alpha_n]$ be a set of logical clauses and let $B = [\beta_1, \cdots, \beta_n]$ the corresponding positive weights. Let $\theta_j$ be a possible assignment of constants (from the domain $D$) to the variables (e.g. X, Y) of the clause $\alpha_i$, that is, a substitution. Let $\alpha_i \theta_j$ the grounded clause where all variables in $\theta_j$ have been replaced by their corresponding constants. Finally, let $\mathbb{1}(\omega, \alpha_i \theta_i)$ be an indicator function, evaluating to 1 if the ground clause is *True* in $\omega$, 0 otherwise.

The probabilistic semantics of Markov Logic is the distribution (with $Z$ the normalization constant):

$$p(\omega) = \frac{1}{Z} \exp \Big( \sum_i \beta_i \sum_j \mathbb{1}(\omega, \alpha_i \theta_j) \Big)$$

Intuitively, in MLNs, a world is more probable if it makes many of its ground instances *True*. Notice that MLNs are usually defined on first-order clause theories, with variables and domains. We will further investigate this issue in Section 3.

## 2.2. Implications for NeSy

The distinction between proof vs models and inference rules vs constraints, turns out to be fundamental for neurosymbolic systems as well.

In neurosymbolic AI, weighted clauses are not used to construct a probabilistic graphical model, but they are likewise used to construct neural models. More specifically, NeSy systems that exploit a proof-theoretic approach use the proofs to build the *architecture* of the neural net. On the other side of the spectrum, NeSy systems that exploit a model-theoretic approach use the constraints to build a *loss function* for the neural net.

Both choices are extremely natural. Proof trees capture the structure of the inference process in a graphical representation. Therefore, they can be used as the structure of the neural network computation, which corresponds to their architecture. On the other hand, the desired behaviour of the variables is expressed in terms of constraints. Loss functions are the de facto standard to enforce desired behaviours on the output variables of a neural network.

First, we survey proof-based NeSy models, which use theorem proving for logical inference and proofs to template the neural architecture. In particular, when proving a specific query atom, they keep track of all the used rules in a proof tree, such as the one shown in Example 1. Weights on facts and rules are then

used to label leaves or edges of the tree, respectively, while real valued activation functions are used to label the AND and OR nodes. The result is a computational graph that can be executed (or evaluated) bottom-up, starting from the leaves up to the root. Generally speaking, the output of the computational graph is a *score* for the query atom. Different semantics can be exploited in building the computational graph, ranging from relaxations of truth values (such as in fuzzy logic) to probabilities (see Section 4). The connection between the proof tree and the neural network suggests schemes for learning the parameters of these models. Indeed, the obtained computational graph is always differentiable. Thus, given a set of atoms that are known to be *True* (resp. *False*), one can maximize (resp. minimize) their score using the corresponding computational graphs. Inference in these models is turned into *evaluation* of the computational graph. The direction of the rules indicates the direction of the evaluation, in the same way as it indicates the direction of inference in logic programming. Among this category are systems based on Prolog or Datalog, such as TensorLog [18], Neural Theorem Provers (NTPs) [102], NLProlog [131], DeepProbLog [72], NLog [121] and DiffLog [112]. Lifted Relational Neural Networks (LRNNs) [116] and $\partial$ILP [39] are other examples of non-probabilistic directed models, where weighted definite clauses are compiled into a neural network architecture in a forward chaining fashion. The systems that imitate logical reasoning with tensor calculus, Neural Logic Programming (NeuralLP) [137] and Neural Logic Machines (NLM) [35], are likewise instances of directed logic. An example of a proof-based NeSy model is given in Example 5.

---

EXAMPLE 5: KNOWLEDGE-BASED ARTIFICIAL NEURAL NETWORKS

Knowledge-Based Artificial Neural Networks (KBANN)[119] is the first method to use definite clausal logic and theorem proving to template the architecture of a neural network.
KBANN turns a program into a neural network in several steps:

1. KBANN starts from a definite clause program and a set of queries.
2. The program is turned into an AND-OR tree using the proofs for the queries.
3. The AND-OR tree is turned into a neural network with a similar structure. Nodes are divided into layers. The weights and the biases are set such that evaluating the network returns the same outcome of querying the program.
4. New hidden units are added. Hidden units play the role of unknown rules that need to be learned. They are initialized with zero weights; i.e. they are inactive.
5. New links are added from each layer to the next one, obtaining the final neural network.

An example of this process is shown in Figure 3. KBANN needs some restrictions over the kind of rules. In particular, the rules are assumed to
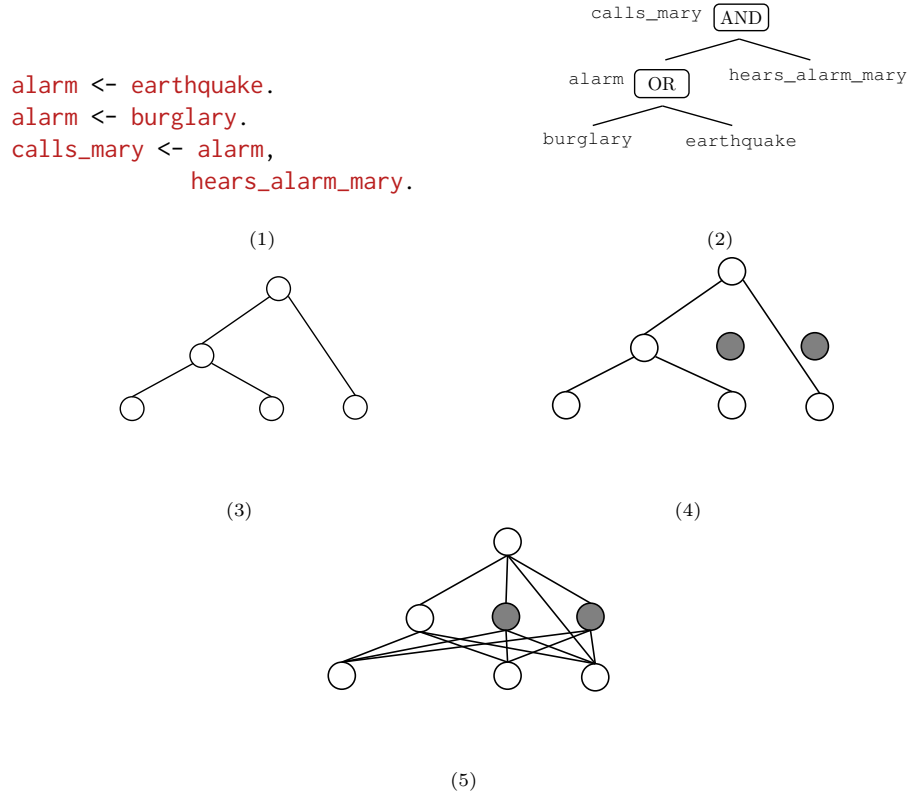
```
alarm <- earthquake.
alarm <- burglary.
calls_mary <- alarm,
          hears_alarm_mary.
```

(1)

calls_mary [AND]

alarm [OR]    hears_alarm_mary

burglary    earthquake

(2)

(3)

(4)

(5)

Figure 3: Knowledge-Based Artificial Neural Network. Network creation process. (1) the initial logic program; (2) the AND-OR tree for the query *calls_mary*; (3) mapping the tree into a neural network; (4) adding hidden neurons, (5) adding interlayer connections.

be conjunctive, non-recursive, and variable-free (or propositional). Many of these restrictions are removed by more recent systems.

We now survey the second class of NeSy systems, the model-based ones. These systems use logic to define a loss function (usually a regularization term) for neural networks. The networks compute scores for the set of atoms that correspond to the output neurons. At each training step, the logic-based loss function determines the degree to which the assigned scores violate the logical theory and uses this to determine the penalty. Logical inference is turned into a learning problem (i.e. "learning to satisfy") and it is usually cast in a variational optimization scheme.[1] As a consequence, in constraint-based models, the neural network has to solve two tasks at the same time: solving a subsymbolic learning

---

[1]This is reminiscent of the variational approach to probabilistic inference in probabilistic graphical models and constitutes a further parallel between the fields.

problem (e.g. perception) as well as approximating the logical inference process [75]. A large group of NeSy approaches, including Semantic Based Regularization (SBR) [33], Logic Tensor Networks (LTN) [5], Semantic Loss (SL) [133] and DL2 [40], exploits logical knowledge as a soft regularization constraint that favours solutions that satisfy the logical constraints. SBR and LTN compute atom (fuzzy) truth assignments as the output of the neural network and translate the provided logical formulas into a real valued regularization loss term using fuzzy logic. SL uses marginal probabilities of the target atoms to define the regularization term and relies on arithmetic circuits [24] to evaluate it efficiently, as detailed in Example 6. DL2 defines a numerical loss providing no specific fuzzy or probabilistic semantics, which allows for including numerical variables in the formulas (e.g. by using a logical term $x > 1.5$). Another group of approaches, including Neural Markov Logic Networks (NMLN) [78] and Relational Neural Machines (RNM) [76] extend MLNs, allowing factors of exponential distributions to be implemented as neural architectures. Finally, [103, 32] compute ground atoms scores as dot products between relation and entity embeddings; implication rules are then translated into a logical loss through a continuous relaxation of the implication operator.

---

EXAMPLE 6: SEMANTIC LOSS

The Semantic Loss [133] is an example of an undirected model where (probabilistic) logic is exploited as a *regularization* term in training a neural model. Let $p = [p_1, \ldots, p_n]$ be a vector of probabilities for a list of propositional variables $X = [X_1, \ldots, X_n]$. In particular, $p_i$ denotes the probability of variable $X_i$ being *True* and corresponds to a single output of a neural net having $n$ outputs. Let $\alpha$ be a logic sentence defined over $X$.
Then, the *semantic loss* between $\alpha$ and $p$ is:

$$Loss(\alpha, p) \propto -\log \sum_{x \models \alpha} \prod_{i:x \models X_i} p_i \prod_{i:x \models \neg X_i} (1 - p_i).$$

The authors provide the intuition behind this loss:

> *The semantic loss is proportional to the negative logarithm of the probability of generating a state that satisfies the constraint when sampling values according to p.*

Suppose you want to solve a multi-class classification task (example adapted from [133]), where each input example must be assigned to a single class. Then, one would like to enforce *mutual exclusivity* among the classes. This can be easily done on supervised examples, by coupling a softmax activation layer with a cross entropy loss. However, there is no standard way to impose this constraint for unlabeled data, which can be useful in a semi-supervised setting.
The solution provided by the Semantic Loss framework is to encode mutual

---

exclusivity into the propositional constraint $\beta$:

$$\beta = (X_1 \wedge \neg X_2 \wedge \neg X_3) \vee (\neg X_1 \wedge X_2 \wedge \neg X_3) \vee (\neg X_1 \wedge \neg X_2 \wedge X_3)$$

Consider a neural network classifier with three outputs $p = [p_1, p_2, p_3]$. Then, for each input example (whether labeled or unlabeled), we can build the semantic loss term:

$$L(\beta, p) = p_1(1 - p_2)(1 - p_3) + (1 - p_1)p_2(1 - p_3) + (1 - p_1)(1 - p_2)p_3$$

It can be summed up to the standard cross-entropy term for the labeled examples.

Unlike for directed methods such as KBANN (Example 5) and TensorLog, the logic is turned into a loss-function that is used during training. The function constrains the underlying probabilities, but there are no directed or causal relationships among them. Moreover, during inference only the probabilities $p$ are used while the logic formula $\beta$ is not used anymore. On the contrary, in KBANN, the logic is compiled into the architecture of the network and, therefore, it it is also exploited at evaluation time.

To conclude, let us stress a key difference between the two classes of NeSy systems w.r.t. *the way they incorporate the knowledge expressed in the logical clauses*. Proof-based, directed models use logic to define the architecture of a neural symbolic network. Thus, logic is part of the inference of the model and acts as a structural constraint. The designer has full control of where and how the logic is used inside the network. Thus, logical knowledge can easily be extended or modified at test-time, without the need to retrain, leading to a high degree of modularity and out-of-distribution generalization [82]. On the other hand, when logic is only encoded in an objective function, the neural net learns to (approximately) satisfy it. Therefore, the knowledge is only latently encoded in the weights of the network, which leads to a loss of control and interpretability. However, the latter techniques are often much more scalable, especially at inference time. The balance between control and interpretability, on the one hand, and scalability, on the other hand, is an open and important research question in the NeSy community.

## 3. Logic - Syntax

In Section 2, we have introduced clausal logic, without paying much attention to the structure of the atoms and literals. This structure and its consequences for StarAI and NeSy models are the topic of the present section. Consider the following example:

Here, the literals do not possess any internal structure. They are *propositions*, which are atoms that we can only assign the value $True$ or $False$. We say that we are working in *propositional logic*.

This contrasts with *first-order* logic in which the literals take the form $p(t_1, ..., t_m)$, with $p$ a predicate of arity $m$ and the $t_i$ terms, that is, constants, variables, or structured terms of the form $f(t_1, ..., t_q)$, where $f$ is a functor and the $t_i$ are again terms. Intuitively, constants represent objects or entities, functors represent functions, variables make abstraction of specific objects, and predicates specify relationships amongst objects. The subset of first order logic where there are no functors is called *relational logic*.

It is interesting to understand the connection between propositional, relational and first-order logic. To this end, we introduce the concept of grounding. When an expression (i.e, clause, atom or term) does not contain any variable it is called *ground*. A substitution $\theta$ is an expression of the form $\{X_1/c_1, ..., X_k/c_k\}$ with the $X_i$ different variables, the $c_i$ terms. Applying a substitution $\theta$ to an expression $e$ (term, atom or clause) yields the instantiated expression $e\theta$ where all variables $X_i$ in $e$ have been simultaneously replaced by their corresponding terms $c_i$ in $\theta$. We can take for instance the clause $mortal(X) \leftarrow human(X)$ and apply the substitution $\{X/socrates\}$ to yield $mortal(socrates) \leftarrow human(socrates)$. Grounding is the process whereby all possible substitutions that ground the clauses are applied. Notice that grounding a first order logical theory may result in an infinite set of ground clauses (when there are functors), and a polynomially larger set of clauses (when working with finite domains).

Finite domains are the focus in both StarAI and NeSy. In such domains, any problem expressed in first-order logic can be equivalently expressed in relational logic and any problem expressed in relational logic can likewise be expressed in propositional logic by grounding out the clauses [94, 41].

## 3.1. Implications for StarAI

StarAI typically focus on first order logic [31, 106, 100]. In Section 2, we have seen how StarAI models can be easily interpreted in terms of probabilistic graphical models (PGM). Here, we want to show that FOL is a powerful tool for building such models.

> FOL allows for knowledge in the form of logical rules to be interpreted as a template for defining the graphical models. Grounding the theory then corresponds to *unrolling* the template. At the same time, first order logic has also an important statistical and learning advantage: a FOL rule leads to parameter sharing in the model as the parameters of a single FOL rule are tied to all its groundings. Parameter sharing compresses the representation of the corresponding probabilistic model, resulting in more efficient learning and better generalization.

These properties are reminiscent of plate notation for probabilistic graphical models, bringing logical reasoning into the picture [97]. In Example 4, we have used only two first order rules but we obtained a larger graphical model with six factors (see Figure 2) by grounding (i.e. unrolling) the rules over the domain. All factors corresponding to the same rule share the same weight.

## 3.2. Implications for NeSy

NeSy exploits the internal structure of literals, resulting in many relational and first-order systems. System exploiting propositional logic are Semantic Loss (SL) [133] and DL2 [40]. Relational logic-based systems are DiffLog [112], $\theta$ILP [39], Lifted Relational Neural Networks (LRNN) [116], Neural Theorem Provers (NTP) [102] and NeurASP [138]. Finally, many systems are based on first-order logic or first-order logic programs, such as DeepProbLog [72], NLog [121], NLProlog [131], DeepStochLog [132], Logic Tensor Networks [5], Semantic Based Regularization [33], Relational Neural Machines [76] and Logical Neural Networks [101].

The focus in NeSy on structured terms is strongly related to that in StarAI and plays a fundamental role in NeSy. In fact, grounding a relational or first-order theory can often be seen as unrolling either the architecture (e.g., DeepStochLog[132], LRNN [116]) or the loss function (e.g., SBR [33], LTN [5]) of the corresponding neural model. Unrolling fixed modules over multiple elements of a complex data structure is fundamental to neural networks on sequences (recurrent nets, RNN), trees (recursive nets, RvNN) and graphs (graph nets, GNN). NeSy can be regarded as unrolling more complex logical structures, with similar benefits in terms of model capacity, modularization and generalization, and strong control due to the formal semantics.

16

Moreover, first-order NeSy models can explicitly deal with how subsymbolic data (e.g. images or audio) are fed to the neural components of the system.In fact, NeSy systems often use subsymbolic data samples as elements of the domain of discourse. For example, the element *mary* can be used to refer to an image, e.g. $mary =$ 🧑. Feeding such samples as input to a neural network can then be naturally encoded as grounding a predicate over the domain of interest. When the internal structure of the literals is absent, as in SL, this mapping must be handled outside the logical framework.

While both relational logic and first-order logic have their advantages, there is a noteworthy distinction in the latter. First-order logic allows representing real valued functions through the use of functors. For example, segmentation can be modeled as a functor returning the bounding box of an object inside an image, e.g. *location(mary,image)* [115]. Therefore, FOL-based systems can address regression tasks, diverging from the conventional classification tasks associated with relational logic systems.

## 4. Logic - Semantics

### 4.1. Model-theoretic semantics

The semantics of logical, probabilistic logical and neurosymbolic systems is defined in terms of a model theoretic semantics. In the present section, we will restrict our attention to Herbrand interpretations and models as is usual in logic programming and statistical relational AI (see Section 2).

We can distinguish three different levels of semantics, which are also closely tied to the used syntax of the underlying logic.

First, when the logical theory consists of definite clauses only, the semantics is given by the least Herbrand model. The least Herbrand model of a definite clause theory is unique and it is the smallest w.r.t. set inclusion. It contains all ground facts (from the Herbrand domain) that are logically entailed by the theory. For instance, considering the facts $a$ and $b$ and the rules $d \leftarrow a, b$ and $b \leftarrow c$ would give the least Herbrand model $\{a, b, d\}$.

Second, when the logical theory can contain any set of clauses, the semantics is given by the set of all possible Herbrand models. For instance, considering the clause $a \vee b$ yields the models $\{a\}, \{b\}$ and $\{a, b\}$. So there is not necessarily a unique model, not even when considering only minimal models, where we have $\{a\}, \{b\}$.

Third, while Horn-clauses are the basis for "pure" Prolog and logic programs, there exist several extensions to this formalism to accommodate negated literals in the condition part of rules or disjunction in the head. A popular framework in this regard is answer set programming (ASP). In ASP the clause $a \vee b$ could be represented by two clauses $a \leftarrow \neg b$ and $b \leftarrow \neg a$ which would have two stable models $\{a\}$ and $\{b\}$.

### 4.2. Fuzzy semantics

The previous three levels of semantics are based on Boolean models, i.e. models where each atom is either present (i.e. *True*) or absent (i.e. *False*).

17

Differently, *fuzzy logic*, and in particular t-norm fuzzy logic, assigns a truth value to atoms in the *continuous* real interval $[0, 1]$. Logical operators are then turned into real-valued functions, mathematically grounded in the t-norm theory. A *t-norm* $t(x, y)$ is a real function $t : [0, 1] \times [0, 1] \rightarrow [0, 1]$ that models the logical AND and from which the other operators can be derived. Table 2 shows well-known t-norms and the functions corresponding to their connectives. A fuzzy logic formula is mapped to a real valued function of its input atoms, as we show in Example 9. Fuzzy logic generalizes Boolean logic to continuous values. All the different t-norms are coherent with Boolean logic in the endpoints of the interval $[0, 1]$, which correspond to *completely true* and *completely false* values. The concept of model in fuzzy logic can be easily recovered from an extension of the model-theoretic semantics of the Boolean logic. Any *fuzzy* interpretation is a model of a formula if the formula evaluates to 1.

EXAMPLE 9: FUZZY LOGIC

Let us consider the following propositions: *alarm*, *burglary* and *earthquake*. Defining a fuzzy semantics for this language requires one to assign truth degrees to each of the propositions and selecting a particular t-norm to implement the connectives.

Let us consider the Łukasiewicz t-norm and the following interpretation of the language:

$$\mathcal{I} = \{\mathsf{alarm} = 0.7,$$
$$\mathsf{burglary} = 0.6,$$
$$\mathsf{earthquake} = 0.3\}$$

$$t_\vee(x, y) = \min(1, x + y)$$
$$t_\rightarrow(x, y) = \min(1, 1 - x + y)$$
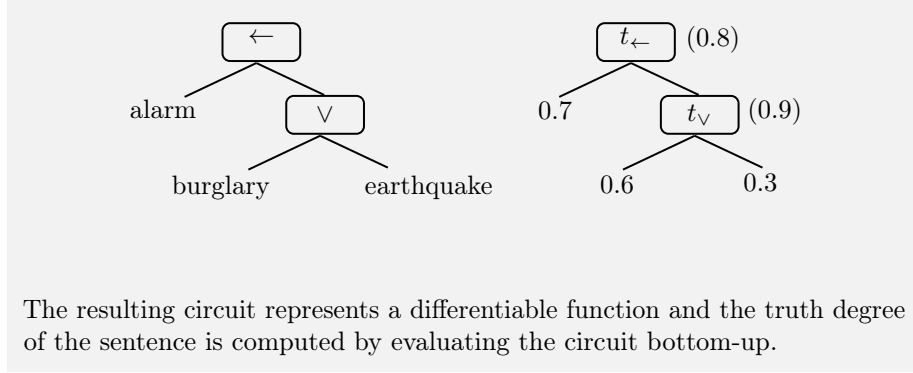
Once we have defined the semantics of the language, we can evaluate logic sentences, e.g.:

$$\mathsf{alarm} \leftarrow (\mathsf{burglary} \vee \mathsf{earthquake}) =$$
$$\min(1, 1 - \min(1, \mathsf{burglary} + \mathsf{earthquake}) + \mathsf{alarm}) = 0.8$$

This evaluation can be performed automatically by parsing the logical sentence in the corresponding *expression tree* and then compiling the expression tree using the corresponding t-norm operation:

| | Product | Łukasiewicz | Gödel |
|---|---|---|---|
| $x \wedge y$ | $x \cdot y$ | $\max(0, x + y - 1)$ | $\min(x, y)$ |
| $x \vee y$ | $x + y - x \cdot y$ | $\min(1, x + y)$ | $\max(x, y)$ |
| $\neg x$ | $1 - x$ | $1 - x$ | $1 - x$ |
| $x \Rightarrow y \;\; (x > y)$ | $y/x$ | $\min(1, 1 - x + y)$ | $y$ |

Table 2: Logical connectives on the inputs $x, y$ when using the fundamental t-norms.



The resulting circuit represents a differentiable function and the truth degree of the sentence is computed by evaluating the circuit bottom-up.

### 4.3. Implications for StarAI

Statistical Relational AI has extended the previous semantics by defining probability distributions $p(\omega)$ over models, or *possible worlds*[2].

The goal is to reason about the uncertainty of logical statements. In particular, the probability that a certain formula $\alpha$ holds is computed as the sum of the probabilities of the possible worlds that are models of $\alpha$ (i.e. where $\alpha$ is True):

$$p(\alpha) = \sum_{\omega \models \alpha} p(\omega) \tag{1}$$

This is an instance of the Weighted Model Counting (WMC) problem. In fact, we are counting how many worlds are models of $\alpha$ and we are weighting each of them by its probability according to the distribution $p(\omega)$.

EXAMPLE 10: PROBABILISTIC LOGIC

Let us consider the following set of propositions $B = burglary$, $E = earthquake$, $J = hears\_alarm\_john$ and $M = hears\_alarm\_mary$. In probabilistic logic, a probability distribution over all the possible worlds is defined. For example, Table 3 represents a valid distribution.

---

[2]In this paper, we use the distribution semantics as representative of the probabilistic approach to logic. While this is the most common solution in StarAI, many other solutions exist [85, 52], whose description is out of the scope of the current survey. A detailed overview of the different flavours of formal reasoning about uncertainty can be found in [53].

| B | E | J | M | $p(\omega)$ |
|---|---|---|---|---|
| F | F | F | F | 0.2394 |
| F | F | F | T | 0.1026 |
| F | F | T | F | 0.3591 |
| F | F | T | T | 0.1539 |
| F | T | F | F | 0.0126 |
| F | T | F | T | 0.0054 |
| F | T | T | F | 0.0189 |
| F | T | T | T | 0.0081 |
| T | F | F | F | 0.0266 |
| T | F | F | T | 0.0114 |
| T | F | T | F | 0.0399 |
| T | F | T | T | 0.0171 |
| T | T | F | F | 0.0014 * |
| T | T | F | T | 0.0006 * |
| T | T | T | F | 0.0021 * |
| T | T | T | T | 0.0009 * |

Table 3: A distribution over possible worlds for the four propositional variables *burglary* (B), *earthquake* (E), *hears_alarm_john* (J) and *hears_alarm_mary* (M). The ∗ indicates those worlds where *burglary* ∧ *earthquake* is *True*.

Suppose we want to compute the probability of the formula *burglary* ∧ *earthquake*. This is done by summing up the probabilities of all the worlds where both *burglary* and *earthquake* are *True* (indicated by a ∗ in Table 3).

The StarAI community has provided several formalisms to define such probability distributions over possible worlds using labeled logic theories. Probabilistic Logic Programs (cf. Example 3) and Markov logic networks (cf. Example 4) are two prototypical frameworks. For example, the distribution in Table 3 is the one modeled by the ProbLog program in Example 3.

It is interesting to compare Markov Logic (Example 4) to ProbLog (Example 3) in terms of their model-theoretic semantics. Markov Logic is defined as a set of weighted full clauses, i.e. as an unnormalized probability distribution over full clausal theories. This means that, given any subset of the theory, there can be many possible models. For instance, the theory $a \lor b$, has three possible models. To obtain a probability distribution over models, Markov Logic needs to distribute the probability mass over its models. To do this, the maximum entropy principle is used, which results in equal distributions of the probability mass. Conversely, ProbLog defines a probability distribution over definite clause theories, each obtained as subsets of the provided probabilistic facts. However, since each of these theories has a unique least Herbrand model, the probability mass corresponding to the selected facts is assigned to the corresponding unique Herbrand model. This means that when working with definite clauses only, there is no need to distribute the probability mass to multiple models and, therefore,
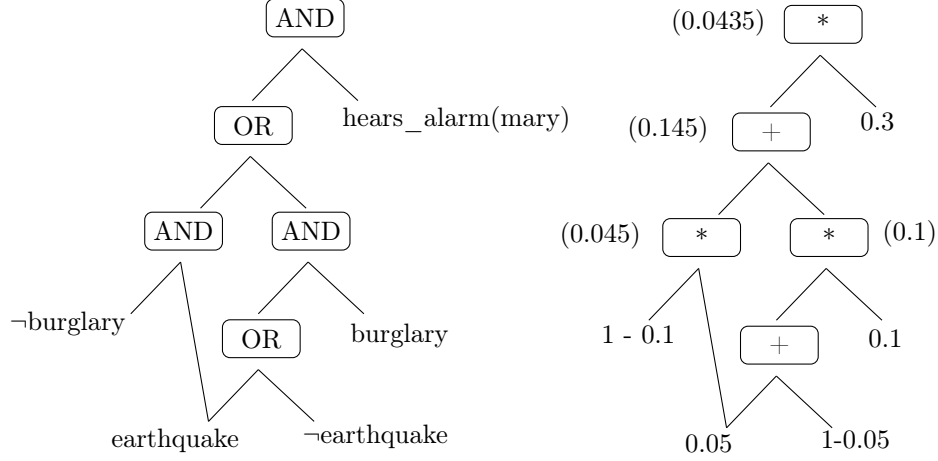
Figure 4: dDNNF (left) and arithmetic circuit (right) corresponding to the ProbLog program in Example 3

no extra assumptions such as maximum entropy are necessary.

Probabilistic inference (i.e. weighted model counting) is generally intractable. That is why, in StarAI, techniques such as *knowledge compilation* (KC) [25] are used. Knowledge compilation transforms a logical formula $\alpha$ into a new representation in an offline step, which can be computationally expensive. Using this new representation a particular set of queries can be answered efficiently (i.e. in poly-time in the size of the new representation). From a probabilistic point of view, this translation solves the disjoint-sum problem, which states that one cannot simply sum up the probability of two disjuncts but also has to subtract the probability of the intersection. After the translation, the probabilities of any conjunction and of any disjunction can be simply computed by multiplying, resp. summing up, the probabilities of their operands. Thus a logical formula $\alpha$ can be compiled into an arithmetic circuit $ac(\alpha)$. The weighted model count of the query formula can then simply be computed by evaluating the corresponding arithmetic circuit bottom up; i.e. $p(\alpha) = ac(\alpha)$.

EXAMPLE 11: KNOWLEDGE COMPILATION

Let us consider the ProbLog program in Example 3 and the corresponding tabular representation in Table 3. Let us consider the query $q = calls(mary)$. Now we can use Equation 1 to compute the probability $p(q)$. To do this, we iterate over the table and we sum all the probabilities of the worlds where $calls(mary)$ is True, which we know from Example 1 are those where either $burglary = T$ or $earthquake = T$ and where $hears\_alarm(mary) = T$. This yields $p(q) = 0.0435$. This method would require us to iterate over $2^N$ terms

(where $N$ is the number of probabilistic facts).

Knowledge compilation compiles $\alpha$ into some normal form that is logically equivalent. In Figure 4, the target representation is a decomposable, deterministic negative normal form (d-DNNF) [23], for which weighted model counting is poly-time in the size of the formula. Decomposability means that, for every conjunction, the two conjuncts do not share any variables. Deterministic means that, for every disjunction, the two disjuncts are mutually exclusive, i.e., only one of the disjuncts can be true at the same time. The formula in d-DNNF can then be straightforwardly turned into an arithmetic circuit by substituting AND nodes with multiplication and OR nodes by summation. In Figure 4, we show the d-DNNF and the arithmetic circuit of the distribution defined by the ProbLog program in Example 3. The bottom-up evaluation of this arithmetic circuit computes the correct marginal probability $p(\alpha)$ much more efficiently than the naive iterative sum that we have computed before.

Even though probabilistic Boolean logic is the most common choice in StarAI, some approaches use probabilistic fuzzy logic. The most prominent approach is Probabilistic Soft Logic (PSL) [3], illustrated in Example 12. Similarly to Markov logic networks, Probabilistic Soft Logic (PSL) defines log linear models where features are represented by ground clauses. However, PSL uses a fuzzy semantics of the logical theory. Therefore, atoms are mapped to real valued variables and ground clauses to real valued factors.

Example 12: Probabilistic Soft Logic

Let us consider the logical rule $\alpha = smokes(X) \leftarrow stress(X)$ with weight $\beta$. As we have seen in Example 4, Markov Logic translates the formula into a discrete factor by using the indicator functions $\mathbb{1}(\omega, \alpha\theta)$:

$$\phi^{MLN}(\omega, \alpha) = \beta\mathbb{1}(\omega, \alpha\{X/mary\}) + \beta\mathbb{1}(\omega, \alpha\{X/john\})$$

Instead of discrete indicator functions, PSL [3] translates the formula into a continuous t-norm based function:

$$t(\omega, \alpha) = \min(1, 1 - stress(X) + smokes(X))$$

and the corresponding potential is then translated into the continuous and differentiable function:

$$\phi^{PSL}(\omega, \alpha) = \beta t(\omega, \alpha\{X/mary\}) + \beta t(\omega, \alpha\{X/john\})$$

Another important task in StarAI is MAP inference. In MAP inference, given the distribution $p(\omega)$, one is interested in finding the interpretation $\omega^\star$ where $p$ is maximal, i.e.

$$\omega^\star = \arg\max_\omega p(\omega) \tag{2}$$

22

When the $\omega$ is a boolean interpretation, i.e. $\omega \in \{0,1\}^n$, like in ProbLog or MLNs, this problem is related to maxSAT, which is NP-hard. However, in PSL, $\omega$ is a fuzzy interpretation, i.e. $\omega \in [0,1]^n$ and $p(\omega) \propto \exp\left(\sum_i \beta_i \phi(\omega, \alpha_i)\right)$ is a continuous and differentiable function. The MAP inference problem can thus be *approximated* more efficiently than its boolean counterpart using gradient-based techniques.

### 4.4. Implications for NeSy

We have seen that in StarAI, one can turn inference tasks into the evaluation (as in KC) or gradient-based optimization (as in PSL) of a differentiable parametric circuit. The parameters are scalar values (e.g. probabilities or truth degrees) that are attached to basic elements of a logical theory (facts or clauses).

A natural way of carrying over the StarAI approach to NeSy is the reparameterization method. Reparameterization substitutes the scalar values assigned to facts or formulas with the output of a neural network. One can interpret this substitution in terms of a different parameterization of the original model. Many probabilistic methods parameterize the underlying distribution in terms of neural components. In particular, as we show in Example 13, DeepProbLog exploits neural predicates to compute the probabilities of probabilistic facts as the output of neural computations over vectorial representations of the constants, which is similar to SL in the propositional counterpart (see Example 6). NeurASP also inherits the concept of a neural predicate from DeepProbLog.

EXAMPLE 13: PROBABILISTIC SEMANTICS REPARAMETERIZATION IN DEEP-PROBLOG

DeepProbLog [72] is a neural extension of the probabilistic logic programming language ProbLog. DeepProbLog allows images or other subsymbolic representations as terms of the program.
Let us consider a possible neural extension of the program in Example 3. We could extend the predicate $calls(X)$ with two extra inputs, i.e. $calls(B, E, X)$. $B$ is supposed to contain an image of a security camera, while $E$ is supposed to contain the time-series of a seismic sensor. We would like to answer queries like $calls(\boxed{\text{▯}}, \boxed{\text{⏸}}, mary)$, i.e. what is the probability that $mary$ calls, given that the security camera has captured the image $\boxed{\text{▯}}$ and the sensor the signal $\boxed{\text{⏸}}$ .
DeepProbLog can answer this query using the following program:

```
nn(nn_burglary, [B]) :: burglary(B).
nn(nn_earthquake, [E]) :: earthquake(E).
0.3::hears_alarm(mary).
0.6::hears_alarm(john).
alarm(B,_) <- burglary(B).
alarm(_,E) <- earthquake(E).
calls(B,E, X) <- alarm(B,E), hears_alarm(X).
```

Here, the program has been extended in two ways. First, new arguments (i.e. $B$ and $E$) have been introduced in order to deal with the subsymbolic inputs. Second, the probabilistic facts *burglary* and *earthquake* have been turned into *neural predicates*. Neural predicates are special probabilistic facts that are annotated by neural networks instead of by scalar probabilities.

Inference in DeepProbLog mimics that of ProbLog. Given the query and the program, knowledge compilation is used to build the arithmetic circuit in Figure 5.

Since the program is structurally identical to the purely symbolic one in Example 11, the arithmetic circuit is exactly the same. The only only difference is that some leaves of the tree (i.e. capturing probabilities of facts) can now also be neural networks.

Given a set of queries that are *True*, i.e.:

$$\mathcal{D} = \{calls(\blacksquare, \blacksquare, mary),$$
$$calls(\blacksquare, \blacksquare, john),$$
$$calls(\blacksquare, \blacksquare, mary), ...\},$$

we can train the parameters $\theta$ of the DeepProbLog program (both neural networks and scalar probabilities) by maximizing the log-likelihood of the training queries using gradient descent:

$$\max_{\theta} \sum_{q \in \mathcal{D}} \log p(q)$$

Similarly to DeepProbLog, NMLNs and RNMs use neural networks to parameterize the factors (or the weights) of a Markov Logic Network. [103] computes marginal probabilities as logistic functions over similarity measures between embeddings of entities and relations. An alternative solution to exploit a probabilistic semantics is to use knowledge graphs (see also Appendix A) to define probabilistic priors to neural network predictions, as done in [118].

SBR[33] and LTN[5] reparametrize fuzzy atoms using neural networks that take as inputs the feature representation of the constants and return the corresponding truth value, as shown in Example 14. Logical rules are then relaxed into soft constraints using fuzzy logic. Many other systems exploit fuzzy logic to inject knowledge into neural models [48, 68]. These methods can be regarded as variants of a unique conceptual framework as the differences are often minor and in the implementation details.

EXAMPLE 14: SEMANTIC-BASED REGULARIZATION

Semantic-Based Regularization (SBR) [33] is an example of an undirected model where fuzzy logic is exploited as a *regularization* term when training a neural model.

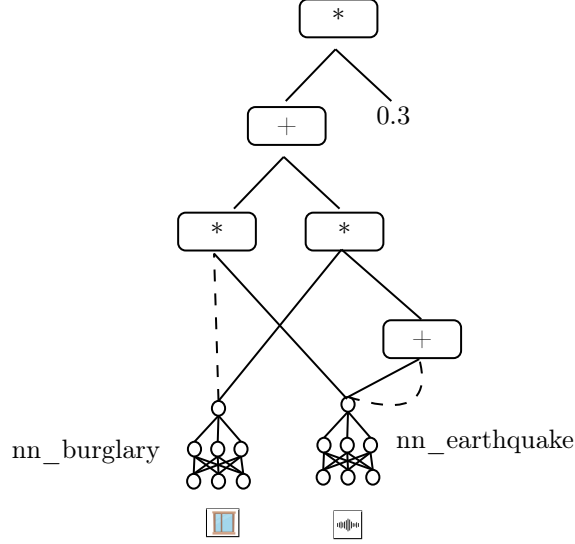Let us consider a possible grounding for the rule in Example 12:

Figure 5: A neural reparametrization of the arithmetic circuit in Example 11 as done by DeepProbLog (cf. Example 13). Dashed lines indicate a negative output, i.e 1 - x. We use a different notation for negation than in Figure 4 to stress that both leaves are parameterized by the same neural network.

$$\text{smokes(mary)} \leftarrow \text{stress(mary)}$$

For each grounded rule $r$, SBR builds a regularization loss term $L(r)$ in the following way. First, it maps each constant $c$ (e.g. *mary*) to a set of (perceptual) features $x_c$ (e.g. a tensor of pixel intensities $x_{\mathsf{mary}}$). Each relation $r$ (e.g. *smokes, stress*) is then mapped to a neural network $f_r(x)$, where $x$ is the tensor of features of the input constants and the output is a truth degree in $[0, 1]$. For example, the atom *smokes(mary)* is mapped to the function call $f_{\mathsf{smokes}}(x_{\mathsf{mary}})$. Then, a fuzzy logic t-norm is selected and logic connectives are mapped to the corresponding real valued functions. For example, when the Łukasiewicz t-norm is selected, the implication is mapped to the binary real function $f(x, y) = \min(1, 1 - x + y)$.

For the rule above, the Semantic-Based Regularization loss term is (for the Łukasiewicz t-norm):

$$L^{\text{Ł}}(r) = \min\left(1, 1 - f_{\mathsf{stress}}(x_{\mathsf{mary}}) + f_{\mathsf{smokes}}(x_{\mathsf{mary}})\right)$$

The aim of Semantic-Based Regularization is to use the regularization term together a with classical loss function for supervised learning to learn the functions associated to the relations (here $f_{\mathsf{stress}}$ and $f_{\mathsf{smokes}}$).

It is worth comparing this method with the Semantic Loss (Example 6). Both

methods turn a logic formula (either propositional or first-order) to a real valued function that is used as a regularization term. However, because of the different semantics, these two methods have different properties. On the one hand, SL preserves the original logical semantics, by using probabilistic logic. However, due to the probabilistic assumption, the input formula cannot be compiled directly into a differentiable loss but needs to be first translated, i.e. compiled, into an equivalent deterministic and decomposable formula. While this step is necessary for the probabilistic model to be sound, the size of the resulting formula can be exponential in the size of the grounded theory. On the other hand, in SBR, the formula can be compiled directly into a differentiable loss, whose size is linear in the size of the grounded theory. However, in order to do so, the semantics of logic is altered, by turning it into fuzzy logic.

Fuzzy logic can also be used to relax rules. For example, in LRNN[116], $\partial$ILP[39], DiffLog[112] and the approach of [129], the scores of the proofs are computed using fuzzy logic connectives. The theory t-norms has identifying parameterized (i.e. weighted) classes of t-norms [117, 101] that are very close to standard neural computation patterns (e.g. ReLU or sigmoidal layers). This creates an interesting, still not fully understood, connection between soft logical inference and inference in neural networks. A large class of methods [80, 32, 18, 131] relaxes logical statements numerically, without explicitly defining a specific semantics. Usually, the atoms are assigned scores in $\mathbb{R}$ computed by a neural scoring function over embeddings. Numerical approximations are then applied either to combine these scores according to logical formulas or to aggregate proofs scores. The resulting neural architecture is usually differentiable and, thus, trained end-to-end.

Some NeSy methods, such as PSL, have used mixed probabilistic and fuzzy semantics. In particular, Deep Logic Models (DLM)[77] extend PSL by adding neurally parameterized factors to the Markov field, while [57] uses fuzzy logic to train posterior regularizers for standard deep networks using knowledge distillation [55].

The semantics of computational logic has also been explored and extended along other directions that have also been used within AI, for example, *modal* and *temporal* logics [125]. While their analysis is out of the scope of the paper, it is worth mentioning that also such formalisms have been investigation from a neurosymbolic perspective [29, 30, 51].

## 5. Structure versus Parameter Learning

Learning approaches in StarAI and NeSy are usually distinguished as to whether the structure [64] or the parameters of the model are learned [49, 70]. In structure learning, the learning task is to discover the logical theory, i.e., a set of logical clauses and their corresponding probabilities or weights that reliably explains the examples. What *explaining the examples* exactly means depends

on the learning setting. In discriminative learning, we are interested in learning a theory that explains, or predicts, a specific target relation given background knowledge. In generative learning, there is no specific target relation; instead, we are interested in a theory that explains the interactions between all relations in a dataset. In contrast to structure learning, parameter learning starts with a given logical theory and only learns the corresponding probabilities or weights.

Structure learning is an inherently NP-complete problem of searching for the right combinatorial structure, whereas parameter learning can be achieved with any curve fitting technique, such as gradient descent or least-squares. While parameter learning is, in principle, an easier problem to solve, it comes with a strong dependency on the provided user input. If the provided clauses are of low quality, the resulting model will also be of low quality. Structure learning, on the other hand, is less dependent on the user provided input, but is an inherently more difficult problem.

### 5.1. Implications for StarAI

Both structure and parameter learning are common in StarAI. Structure learning in StarA is an instance of learning by search [83] and is closely connected to program synthesis. The existing techniques are typically extensions of techniques originating in inductive logic programming (ILP) [86, 94], which learn deterministic logical theories, and probabilistic graphical models (PGMs), which learn Bayesian or Markov networks from data. Being an instance of learning by search, the central components of a learning framework are a space of valid structures and a search procedure. In ILP, valid structures are logical theories; for Bayesian networks, valid structures are DAGs capturing their graph structure. The resulting search space is then traversed with generic search procedures.

StarAI structure learning techniques suffer from a combinatorial explosion. That is especially the case with ILP techniques, in which the search space consists of programs containing several clauses. Therefore, it is necessary to limit the search space to make learning tractable. The most common way to do this is to impose a *language bias* – a set of instructions on how to construct the search space, such that it is narrowed down to a subset of the space of all logical theories. Though language bias can make the problem more tractable, it requires special care: too many restrictions might eliminate the target theory, while too few restrictions make the search space too large to traverse. Another strategy is to leverage the compositionality of logic programs: adding an additional clause to a program increases its coverage and does not affect the prediction of examples covered by the initial program. That is, we can learn a single clause at a time instead of simultaneously searching over theories containing multiple clauses.

Learning clauses and their probabilities is usually treated as a two stage process. ILP-based StarAI learning techniques first identify useful (deterministic) clauses, and then learn the corresponding probabilities or weights via parameter learning. Similarly, StarAI methods grounded primarily in PGMs, such as MLNs, search for frequently occurring cliques in data [65], lift them into logical clauses, and then learn the weights or probabilities. Parameter learning techniques are

often also extensions of well known statistical approaches such as least-squares regression [49], gradient descent [70], and expectation maximisation [50].

EXAMPLE 15: STRUCTURE LEARNING WITH PROBFOIL

As an illustration of structure learning techniques, we will focus on Prob-Foil [95]. Assume that we are interested in learning the definition of *grandparent* from akinship data. That is, we are given a set of examples of grandparent relations

```
grandparent(abe,lisa).
grandparent(abe,bart).
grandparent(jacqueline,lisa).
grandparent(jacqueline,bart).
```

and background knowledge containing the following facts:

```
father(homer,lisa).
father(homer,bart).
father(abe,homer).
mother(jacqueline,marge).
mother(marge,bart).
mother(marge,lisa).
```

ProbFoil iteratively searches for a single clause that covers as many examples as possible, until all examples are covered or it adding more clauses does not improve the results. While searching for the best clause, it starts from the most general one, *grandparent(X,Y)*, which effectively states that every pair of people forms a grandparent relationship. Then it gradually specialises the clause by adding literals to the body. For instance, extending *grandparent(X,Y)* with a *mother/2* predicate results in the following clauses

```
grandparent(X,Y) <- mother(X,Y).
grandparent(X,Y) <- mother(X,X).
grandparent(X,Y) <- mother(Y,X).
grandparent(X,Y) <- mother(Y,Y).
```

Extending the initial clause with the *father/2* results in similar clauses. Having the new candidate clauses, ProbFoil scores each candidate by counting how many positive and negative examples are covered. These candidate clauses would not cover any examples and ProbFoil continues to refine the candidates by adding another literal to the body. This would result in clauses of the following form:

```
grandparent(X,Z) <- mother(X,Y), father(Y,Z).
grandparent(X,Z) <- mother(X,Y), mother(Y,Z).
grandparent(X,Z) <- father(X,Y), mother(Y,Z).
...
```

Some of the new candidates will cover only positive examples, such as

```
    grandparent(X,Z) <- mother(X,Y), mother(Y,Z)
```

that covers both examples

```
    grandparent(jacqueline,lisa).
    grandparent(jacqueline,bart).
```

Having found one clause, ProbFoil learns the corresponding probability labels and adds the clause to the theory. ProbFoil then repeats the same procedure, starting with the most general clause, to cover the remaining examples.

*5.2. Implications for NeSy*

While StarAI learning techniques are categorised exclusively as either structure or parameter learning, NeSy learning techniques combine both. We will now discuss four groups of NeSy learning approaches: neurally-guided search, structure learning via parameter learning, program sketching, and implicitly structure learning.

*Neurally guided structure search* [60, 37, 38, 122] is the NeSy paradigm most similar to structure learning in StarAI. It addresses one of the major weaknesses of StarAI structure learning methods - uninformed search over valid theories. Instead, neurally guided search relies on a *recognition model*, typically a neural network, to prioritise parts of the symbolic search space so that the target model can be found faster. Generally speaking, the recognition model predicts the probability of a certain structure, e.g. a predicate or an entire clause, to be a part of the target model. For instance, Deepcoder [6] uses input-output examples to predict the probability of each predicate appearing in the target model. Therefore, Deepcoder turns a systematic search into an informed one by introducing a ranking over predicates in the search space. Likewise, $EC^2$ [37] derives the probability of a program solving the task at hand. Several approaches push this direction further and explore the idea of replacing an explicit symbolic model space with an implicit generative model over symbolic models [90, 74]. For instance, in [90], the authors learn a generative model over grammar rules, conditioned on the examples. Structure learning is then performed by sampling grammar rules from the generative model, according to their probability, and evaluating them symbolically on the provided examples.

These approaches clearly show how symbolic search can be made tractable by introducing various forms of guidance via neural models. These guidance-based approaches reduce, to a large extent, the most important weakness of symbolic structure learning approaches - the generation of many useless clauses or models. On the other hand, these approaches often need large amounts of data for training, sometimes millions of examples [38] even though creating data is relatively easy by enumerating random model structures and sampling examples from them [38].

EXAMPLE 16: NEURALLY-GUIDED STRUCTURE LEARNING

To illustrate neurally-guided search, we use the approach of Zhang et al. [140]. StarAI techniques for structure learning typically perform a systematic search, which results in many useless models being tested. Given $N$ atoms, we can construct $N^l$ clauses of length $l$; this is an enormous space that is difficult to search efficiently.

Zhang et al. sidestep the systematic search by introducing a neural network that chooses which programs to explore next. This search space is made of clauses such that an empty clause is at the top and children are extensions of the empty clause with all possible predicates; their children are further extensions with all individual atoms.

The approach follows a top-down search strategy, exploring shorter clauses before longer ones, with a twist: instead of following a predefined order, the approach uses a neural network to decide which child to expand next. The approach can be viewed as a best-first search with a heuristic function implemented by a neural model. To this end, the network (1) encodes all literals in each clause separately, (2) scores all literals, (3) pools the scores of each literal per candidate, and (4) chooses the best candidate based on their scores. Ordering the search space in this way leads to substantial improvements in computation time, typically several orders of magnitude.

An alternative way to reduce the combinatorial complexity of learning is to learn only a part of the program. This is known as *program sketching:* a user provides an almost complete target model with certain parts being unspecified (known as *holes*). For instance, when learning a model in the form of a (logic) program for sorting numbers or strings, the user might leave the comparison operator unspecified and provide the rest of the program. The learning task is then to fill in the holes. Examples of NeSy systems based on sketching are DeepProbLog and $\partial 4$, which fill in the holes in a (symbolic) program via neural networks.

The advantage of sketching is that it provides a nice interface for NeSy systems, as the holes can be filled either symbolically or neurally. Holes provide a clear interface in terms of inputs and outputs and are agnostic to the specific implementation. The disadvantage of sketching is that the user still needs to know, at least approximatively, the structure of the program. The provided structure, the sketch, acts as a strong bias. Deciding which functionality is left as a hole is a non-trivial issue: as the sketch becomes less strict, the search space becomes larger.

*Structure learning via parameter learning* (Example 17) is arguably the most prominent learning paradigm in NeSy, positioned in between the two StarAI learning paradigms. Structure learning via parameter learning is technically equivalent to parameter learning in that the learning tasks consists of learning the probabilities of a fixed set of clauses. However, in contrast to StarAI in

which the user carefully selects the informative clauses, the clauses are typically enumerated from user-provided templates of predefined complexity. Constructed in this way, the majority of clauses are noisy and erroneous and are of little use. They would receive very low, but non-zero, probabilities. Approaches that follow this learning principle include NTPs [102], $\partial$ILP [39], DeepProbLog[72], NeuralLP [137] and DiffLog [112].

The advantage of structure learning via parameter learning is that it removes the combinatorial search from the learning. However, the number of clauses that needs to be considered is still extremely large, which leads to difficult optimisation problems (cf. [39]). Furthermore, irrelevant clauses are never removed from the model and are thus always considered during inference. This can lead to spurious interactions even when low probabilities are associated to irrelevant clauses: as the number of irrelevant clauses is extremely large, their cumulative effect can be substantial.

EXAMPLE 17: STRUCTURE LEARNING VIA PARAMETER LEARNING

As an illustration of structure learning via parameter learning, we focus on DiffLog [112]. DiffLog expects the candidate clauses to be provided by the user. The user can either provide the rules she knows are useful or construct them by using a clause template and instantiating it [20]..

Given a set of positive examples, DiffLog proceeds by constructing *derivation trees* for each example. Consider the problem of learning the connectivity relation over a graph. The input tuples (background knowledge in StarAI terminology) specify edges in a graph

```
edge(a,b). edge(b,c). edge(b,d). edge(d,e). edge(c,f).
```

The examples indicate the connectivity relations among the nodes in the graph (for simplicity, consider only the following two examples)

```
connected(a,b). connected(a,c).
```

Also assume that the candidate clause set contains the following clauses (with $p_1$ and $p_2$ their weights):

$$p_1\text{::connected(X,Y) <- edge(X,Y).}$$
$$p_2\text{::connected(X,Y) <- edge(X,Z), connected(Z,Y).}$$

Derivation trees are essentially proofs of individual examples that correspond to branches in the SLD-tree [69]. For instance, the example *connected(a,b)* can be proven using the first clause, whereas the example *connected(a,c)* can be proven by chaining the two clauses ($connected(a,c) \leftarrow edge(a,b), connected(b,c)$ and $connected(b,c) \leftarrow edge(b,c)$).

DiffLog uses derivation trees to formulate the learning problem as numerical optimisation over the weights associated with the rules. More precisely, DiffLog defines the probability of deriving an example as the product of the weights associated to the clauses used in the derivation tree of the corresponding example. For instance, DiffLog would formulate the learning

problem for the two examples as follows

$$\min_{p_1,p_2} \quad \underbrace{(1-p_1)}_{\texttt{connected(a,b)}} + \underbrace{(1-p_1 \times p_2)}_{\texttt{connected(a,c)}}.$$

The last group of approaches learns the structure of a program only *implicitly*. For instance, Neural Markov Logic Networks (NMLN) [78], a generalisation of MLNs, extract structural features from relational data. Whereas MLNs define potentials only over cliques defined by the structure (logical formulas) of a model, NMLNs add potentials over *fragments* of data (projected over a subset of constants). NMLNs thus do not necessarily depend on the symbolic structure of the model, be it learned or provided by a user, but can still learn to exploit relational patterns present in data. Moreover, NMLNs can incorporate embeddings of constants. The benefit of this approach is that it removes combinatorial search from learning and performs learning via more scalable gradient-based methods. However, one loses the ability to inspect and interpret the discovered structure. Additionally, to retain tractability, NMLNs limit the size of fragments which imposes limits on the complexity of the discovered relational structure.

## 6. Symbolic vs subsymbolic representations

In neurosymbolic artificial intelligence, approaches can be characterized by the way they represents entities and relationships in two classes: symbolic methods, where entities are represented using symbols such as strings and natural numbers, and subsymbolic methods, where entities are represented using numerical or distributed representations.

Symbolic representations include constants $(an, bob)$, numbers $(4, -3.5)$, variables $(X, Y)$ and structured terms $f(t_1, ..., t_n)$ where $f$ is a functor and the $t_i$ are constants, variables or structured terms. Structured terms are a powerful construct that can represent arbitrary structures over entities, such as relations, lists or trees. subsymbolic AI systems, such as neural networks, require that entities are represented numerically using vectors, matrices or tensors. Throughout this this section, we will call these subsymbolic representations or subsymbols. subsymbolic AI systems usually require that these representations have a fixed size and dimensionality. Exceptions require special architectures and are still the subject of active research (e.g. RNNs for list-like inputs or GCNs [63] for graph-type inputs).

*Comparing representations.* A powerful and elegant mechanism for reasoning with symbols in logic is *unification*. Essentially, it calculates the most general substitution that makes two symbols syntactically equal, if it exists. This does not allow one to compare two different entities, but allows one to find what two structured terms have in common. For instance, the terms $p(a, Y)$ and $p(X, b)$ can be unified using the substitution $\{X = a, Y = b\}$. Conversely,

due to their numerical nature, calculating the similarity between subsymbols is straightforward. Similarity metrics such as the radial-basis function or distance metrics such as the L1 and L2 norm can be used. However, it is not clear when to decide that two subsymbolically represented entities are the same.

*Translating between representations.* Many systems need to translate back and forth between symbolic and subsymbolic representations. In fact, a lot of research on deep learning is devoted to efficiently representing symbols so that neural networks can properly leverage them. A straightforward example is to translate symbols to a subsymbolic representation that can serve as input for a neural network. Generally, these symbols are replaced by a one-hot encoding or by learned embeddings. Note, however, that this does not imply that the system can perform symbolic manipulation on this input. Rather, it serves as an index to a set of learned, latent embeddings. A more interesting example is encoding relations in subsymbolic space. The wide variety of methods [13, 120, 136] developed for this purpose indicates that this is far from a solved problem. Different encodings have different benefits. For example, TransE [13] encodes relations as vector translations from subject to object embeddings. A disadvantage is that symmetric relations are represented by the null vector, and entities in symmetric relations are pushed towards each other. More complex structures are even harder to represent. For example, there is currently a lot of research in how to utilize graph-structured data in neural networks (cf. Appendix A).

Translating from a subsymbolic representation back to a symbolic one happens, for example, at the end of a neural network classifier. Here, a subsymbolic vector needs to be translated to discrete classes. Generally, this happens through the use of a final layer with a soft-max activation function which then models the confidence scores of these classes as a categorical distribution. However, other options are possible. For example, some methods are only interested in the most likely class, and will use an arg-max instead. Alternatively, a Gumbel-softmax activation can be used as a differentiable approximation of sampling from the categorical distribution.

## 6.1. Implications for StarAI and NeSy

In StarAI systems, the input, intermediate and output representations are all using the same symbolic representations. Although there are StarAI systems that can support numerical values, these are still treated as symbols, which is different than a latent, subsymbolic representation. In neural systems, the input and intermediate representations are subsymbolic. The output representation can be either symbolic (e.g. classifiers) or subsymbolic (e.g. auto-encoders, GANs). The most important aspect of neurosymbolic systems is that they combine symbolic and subsymbolic representations. NeSy systems can be categorized by how they do this. We distinguish several approaches.

In the first approach, the inputs are symbolic, but they are translated to subsymbols in a single translation step, after which the intermediate representations used during reasoning are purely subsymbolic. This approach is followed

33

by the majority of NeSy methods. Some examples include Logic Tensor Networks [5], Semantic-based Regularization [33], Neural Logic Machines [35] and TensorLog [18].

---

EXAMPLE 18: LOGIC TENSOR NETWORKS

Logic tensor networks [5] make this translation step explicit. The authors introduce the concept of a *grounding* (not to be confused with the term grounding used in logic). Here, a grounding is a mapping of all symbolic entities onto their subsymbolic counterpart. More formally, the authors define a grounding as a mapping $\mathcal{G}$ where:

- $\mathcal{G}(c) \in \mathbb{R}^n$ for every constant symbol $c$

- $\mathcal{G}(f) \in \mathbb{R}^{n.m} \to \mathbb{R}^n$ for every function $f$ of arity $n$

- $\mathcal{G}(p) \in \mathbb{R}^{n.m} \to [0,1]$ for every predicate $p$ of arity $n$

The grounding of a clause is then performed by combining the aforementioned groundings using a t-norm.

---

In the second approach, intermediate representations are both symbolic and subsymbolic, but not simultaneously. This means that some parts of the reasoning work on the subsymbolic representation, and other parts deal with the symbolic representation, but not at the same time. This is indicative of NeSy methods that implement an interface between the logic and neural aspect. This approach is more natural for systems that originate from a logical framework such as DeepProbLog [72], NeurASP [138]), ABL [22] and NLog [121].

---

EXAMPLE 19: ABL

In ABL [22] there are three components that function in an alternating fashion. There is a perception model, a consistency checking component and an abductive reasoning component. Take for example the task where there are 3 MNIST images that need to be recognized such that the last is the result of applying an operation on the first two (e.g.  +  = ). The structure of the expression is given as background knowledge, but the exact operation (addition) needs to be abduced. First, the perception model classifies the images into pseudo-labels, using the most likely prediction (i.e. arg-max). The abductive reasoning component then tries to abduce a logically consistent hypothesis. For example, if the digits are correctly classified as 3, 5 and 8, the only logically consistent hypothesis is that the operation is an addition. If this is not possible, there is an error in the pseudo-labels. A heuristic function is then used to determine which pseudo-labels are wrong. The reasoning module then searches for logically consistent pseudo-labels. These revised pseudo-labels are then used to retrain the perception model.

---

In the final approach, intermediate representations are considered simul-

taneously as symbolic and subsymbolic by the reasoning mechanism.This is implemented in only a few methods, such as the NTP[102] and the CTP[81].
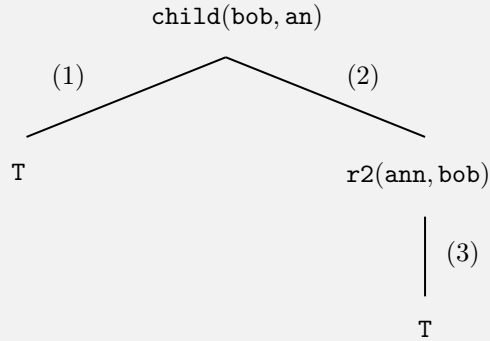
EXAMPLE 20: NEURAL THEOREM PROVER

In the Neural Theorem Prover[102], two entities can be unified if they are similar, and not just if they are identical. As such, the NTP interweaves both symbols and subsymbols during inference. For each symbol $S$, there is a learnable subsymbol $T_S$. Soft-unification happens by applying the normal unification procedure where possible. However, if two symbols $S_1$ and $S_2$ can not be unified, the comparison is assigned a score based on the similarity between $T_{S_1}$ and $T_{S_2}$. The similarity is calculated using a radial basis function $\varphi(||x - y||_2)$.

For example, to unify `mother(an,bob)` and `parent(X,bob)`, soft-unification proceeds as follows:

$$\{\texttt{mother}(\texttt{an}, \texttt{bob}) = \texttt{parent}(\texttt{X}, \texttt{bob})\}$$
$$\Downarrow \quad \varphi(\texttt{mother}, \texttt{parent})$$
$$\{\texttt{an} = \texttt{X}, \texttt{bob} = \texttt{bob}\}$$
$$\Downarrow \quad X = an$$
$$\{\texttt{bob} = \texttt{bob}\}$$
$$\Downarrow$$
$$\{\,\}$$

Soft-unification is not only used to learn which constants and predicates are similar, but can also be used to perform rule learning. By adding new, parameterized rules with unique predicates, soft-unification allows these new predicates to become very similar to other predicates and as such behave as newly introduced rules. For example, consider the program consisting of the fact *mother(an,bob)* and a single parameterized rule $r1(X, Y) \leftarrow r2(Y, X)$. The Neural Theorem Prover can answer the query *child(bob,an)* as follows:

$$\texttt{child}(\texttt{bob}, \texttt{an})$$

(1)           (2)

T                           $\texttt{r2}(\texttt{ann}, \texttt{bob})$

(3)

T

$$
\begin{array}{lll}
(1) & \texttt{child} = \texttt{mother} & \varphi(||T_{child} - T_{mother}||_2) \\
& \texttt{bob} = \texttt{an} & \varphi(||T_{an} - T_{bob}||_2) \\
(2) & \texttt{child} = \texttt{r1} & \varphi(||T_{child} - T_{r1}||_2) \\
(3) & \texttt{child} = \texttt{r2} & \varphi(||T_{r2} - T_{mother}||_2) \\
\end{array}
$$

The figure above shows the two possible derivations the neural theorem prover can make to infer *child(bob, an)*. One the one hand, it can soft-unify with the fact *mother(an, bob)*, where *mother* unifies with *child* and *an* with *bob*. On the other hand, it can use the parameterized rule which encodes an inverse relation. In that case, *mother* unifies with *r1* and *r2* with *child*. If we optimize the subsymbolic embeddings for the latter, this will be equivalent to learning the rule $mother(X, Y) \leftarrow child(Y, X)$. This example also shows that soft-unification potentially adds a lot of different proofs, which can result in computational problems. This problem was solved in later iterations of the system [79].

## 7. Logic vs Probability vs Neural

When two or more paradigms are integrated, examining which of the base paradigms are preserved, and to which extent, tells us a lot about the strengths and weaknesses of the resulting paradigm. It has been argued [98] that when combining different perspectives in one model or framework, such as logic, probabilistic and neural ones, it is desirable to have the original paradigms as a special case.

In this section, we analyze to which extent different models in StarAI and NeSy preserve the three basic paradigms. Intuitively, with preserving we mean to which extent one can exactly replicate the model and inference algorithm of the original paradigm. We will use the capital letters *L, P* and *N* to label systems where the logic, probability and neural paradigms can be recovered in full. We will use lowercase letters (i.e. *l, p* and *n*) when a method only partially recovers these paradigms, i.e. retain some but not all of the features. The absence of a letter means that the paradigm is not considered by an approach.

### 7.1. StarAI: Logic + Probability

Traditionally, StarAI focused on the integration of logic and probability.

***lP:*** The classical knowledge-based model construction approach uses logic only to generate a probabilistic graphical model. Thus the graphical model can be used to define the semantics of the model and also to perform inference. This can make it harder to understand the effects of applying logical inference rules to the model. For instance, in MLNs, the addition of the resolvent of two weighted rules makes it hard to predict the effect on the distribution.

***Lp:*** On the other hand, the opposite holds for probabilistic logic programs (PLPs) and their variants. While the effect of a logical operation is clear, it

is harder to identify and exploit properties such as conditional or contextual independencies, that are needed for efficient probabilistic inference.

## 7.2. NeSy: Logic + Probability + Neural

In NeSy, we consider a third paradigm: neural computation. With neural computation, we refer mainly to the set of models and techniques that allows for exploiting (deep) latent spaces to learn intermediate representations. This includes dealing with perceptual inputs and also dealing directly with embeddings of symbols.

***lN***: Many NeSy approaches focus on the neural aspect (i.e., they originated as a neural method to which logical components have been added). For example, LTNs and SBRs turn the logic into a regularization function to provide a penalty whenever the logical constraints are violated. At test time the logical loss component is dropped and only the network is used to make predictions. Moreover, by using fuzzy logic, these methods do not integrate the probabilistic paradigm.

***Ln***: Another class of NeSy methods does retain the focus on logic. These methods usually expand an existing logical framework into a differentiable version. Examples include LRNNs [116], TensorLog [18], DiffLog [112], $\partial$ILP [39], $\partial$4 [14] and NTPs [102]. The key inference concepts are mapped onto an analogous concept that behaves identically for the edge cases but is continuous and differentiable in non-deterministic cases. As described in the previous sections, many such systems cast logical inference as forward or backward chaining. The focus on logic is clear if one considers that logical inference is performed symbolically to build the network and the semantics is relaxed only in a subsequent stage to learn the parameters. While the architecture mimics the logical reasoning, it is often far from the deep-stacked architecture of neural networks.

***LN***: It is worth mentioning a later iteration of LRNN, where the framework has been extended to allow for tensorial weights on atoms and custom aggregation functions [117]. In that framework, it is shown how specifying logic rules can be regarded as specifying the layers of a deep architecture. This provides a nice and complete integration between forward-chaining logical reasoning and neural networks that is able to implement any existing neural architecture.

***lPN*** and ***LpN*** There are two final classes of methods that start from existing StarAI methods, *lP* and *Lp* respectively, and extend them with primitives that can be interfaced with neural networks and allow for differentiable operations. In the *lPN* class, NeSy methods such as SL, RNMs and NMLNs follow the knowledge-based model construction paradigm. In the *LpN* class, methods such as DeepProbLog and NeurASP extend PLP.

There is usually a trade-off that one must make: systems in the *lN* or *Ln* classes are usually more scalable but *(i)* do not model a probability distribution and *(ii)* often relax the logic. On the contrary, *LpN* or *lPN* systems preserve the original paradigms but at the cost of more complex inference (e.g. they usually resort to exact probabilistic inference).

An aspect that significantly aids in developing a common framework, and analysing its properties, is the development of an intermediate representation language that can serve as a kind of *assembly language* [143]. One such idea concerns performing probabilistic inference by mapping it onto a weighted model counting (WMC) problem. This can then in turn be solved by compiling it into a structure (e.g. an arithmetic circuit) that allows for efficient inference. This has the added benefit that this structure is differentiable, which facilitates the integration between logic based systems and neural networks. StarAI based systems often use this approach.

## 8. Tasks

In this Section, we analyze the learning tasks to which the NeSy models considered in this paper have been applied.

*Distant Supervision.* A classical task in NeSy is to use logic as distant supervision for a learning model. Here, input $X$ is paired with label $y$. However, instead of using a single model to map $X$ to $y$, the input $X$ is firstly mapped to a set of intermediate concepts $C$ by a (set of) neural networks. Then, these concepts are used to compute $y$ in a symbolic way. Logic programs are usually exploited to map the concepts $C$, represented as logical atoms, to the label $y$, which represents the logical query. Therefore, the neural networks are not directly supervised (on $C$) but they are only distantly supervised through the label $y$ and the knowledge contained in the logic program. The intuition is that when the label $y$ is only weakly linked to the input, it is more convenient to break the task in several easier subtasks and then compose them using background knowledge in the form of a logic program. Notice that the logic program is fundamental for the inference. Without the program, the networks will not be able to solve their subtasks, as there is no direct supervision. Moreover, by splitting the task into subtasks, the inference done by the composite system (neural + logic) is far more explainable than a corresponding end-to-end neural network. A classical example is the MNIST addition [72], shown in Example 21. Distant supervision tasks are very common in prototypical systems such as DeepProbLog, DeepStochLog, NLog, NeurASP, SATNet [127]. A downside of such tasks is that, to enable learning of untrained neural subtasks, the logic has to consider all possible combinations of concepts that are compatible with the label $y$, even though only few (or one) are correct. The challenge is to balance the exploration of multiple combinations with a greedy strategy for scaling to larger problems [121, 73, 71]. Other problems falling in this category are scene parsing, image segmentation and semantic image interpretation [34, 2]

---

EXAMPLE 21: MNIST ADDITION

Given the classical MNIST dataset, $\mathcal{D} = \{(x_i, y_i)\}$, with $x_i$ an MNIST image, and $y_i$ its numeric label, the MNIST addition dataset is built by mapping pairs of images to the label representing their sum. In particular,

---

$\mathcal{D}_{\text{add}} = \{(x_i, x_j, z_{ij}) : z_{ij} = y_i + y_j \wedge (x_i, y_i), (x_j, y_j) \in \mathcal{D}\}$. The idea is to learn to classify the digits without direct supervision on their labels, but only using distant supervision about sums of such images. The task is often also coupled to background knowledge of what addition is, e.g. in Prolog syntax:

```
addition(X1, X2, Z) <- digit(X1,Y1), digit(X2,Y2), Z is Y1 + Y2.
```

Such knowledge is used to reason about the (most-likely) pairs `Y1,Y2` that sum to the provided label `Z`. Logic is then used to link the actual outputs of the learning model `Y1,Y2` to the distant supervision `Z`.

*Semi-supervised Classification.* A related class of tasks is semi-supervised classification [15] with knowledge. Here, the starting point is a standard classification task, where a set of inputs $X$ is mapped by a neural model to a set of labels $C$. However, we are also provided with some additional knowledge $y$ related to the labels $C$ of the inputs. This knowledge is often expressed in terms of logical rules and programs. The setting is very similar to distant supervision, where we have three levels: inputs $X$, concepts $C$ and additional labels $y$. However, in this case, we have also access to supervision for some (usually few) concepts $C$. Although this task could be tackled in a purely supervised way by discarding the information contained in $y$, NeSy approaches can improve the predictions of several input patterns using the external knowledge. When the external knowledge is relating concepts $C$ of multiple patterns, the task is called collective classification [109], as one can improve the accuracy on multiple patterns by collectively predicting their classes. A classical example in this setting is document classification in citation networks, cf. Example 22. By treating the information contained in $y$ as extra knowledge, these tasks are often tackled using regularization based systems, like SBR[33], DLM[77], RNM[76] or Semantic Loss[133]. However, logic programs can also be used to simulate a label-passing scheme along the citation network, as done in DeepStochLog [132]. A characteristic of this class of tasks is that the additional information $y$ is often very noisy (e.g. the manifold rule in the citation network is not always valid). While this task is closely related to distant supervision, there is an important difference: in semi-supervised classification, the additional knowledge $y$ is meant to provide an additional signal, which, however, would not suffice in the absence of direct supervision on the concepts $C$.

EXAMPLE 22: DOCUMENT CLASSIFICATION IN CITATION NETWORKS

In document classification in citation networks, we are provided with both labelled and unlabelled scientific papers. A label is often the domain area of the paper (e.g. Machine Learning, Artificial Intelligence, Databases, etc.). However, a network of citations between papers is also provided, linking papers between domains.
The idea of the document classification task is that in many domains, a paper cited by other papers with a certain label is likely to belong to the

same domain. When classifying a document, one has to balance the signal coming from the features of the document (i.e. words) and that coming from neighbors in the citation network to provide a *collective* prediction.

In NeSy systems, this is usually done by coupling the subsymbolic model with a rule of the following type:

```
w:: domain(X,Y) <- cite(X,X1), domain(X1,Y).
```

The rules get a different weight according to the domain to account for the differences between them.

*Knowledge Graph Completion.* Another common task in NeSy is knowledge graph completion (KGC) or link prediction. A knowledge graph (KG) is a pair of $(E, R)$, where $N$ is the set of entities and $R$ the set of edges. In a KG, an edge is a triple $(e_1, r, e_2)$, where $e_1$ and $e_2$ are the head and tail of the edge and $r$ is the relation between them. In a KGC task, the goal is to predict missing edges in the input graph. Link prediction has been one of the key tasks in StarAI [44], and more recently also in NeSy as NeSy allows to merge symbolic reasoning (from StarAI) with the recent geometric deep learning approaches based on Knowledge Graph Embeddings (KGE) [128] and Graph Neural Networks [107]. NeSy systems focusing on this task include NTPs [102], NMLN [78], DLM [77], DiffLog [112], TensorLog [18].

*Generative Tasks.* Most previously mentioned tasks can be described as classification[3]. NeSy has recently also focused on tasks concerned with modeling the input data distribution as accurately as possible. The goal is then to sample new patterns from the learned distribution. The idea behind NeSy generative approaches is that one can learn important features from data using deep generative models (e.g. variational auto-encoders or Markov Chain Monte Carlo methods). Combining symbolic features with logic reasoning can be used to control, stratify and simplify the inference. The generative modeling can either refer to the relational structure, e.g. molecule generation in NMLNs [78], or to the subsymbolic space, e.g. image generation in VAEL [82] or [114].

*Knowledge Induction.* Rather than exploiting symbolic knowledge predictive tasks, one can also induce *symbolic knowledge*. In all previous tasks, symbolic knowledge is provided by the user as part of the input. However, as explored in Section 5, we can still apply several neurosymbolic techniques by learning the symbolic knowledge when this is not the case. The unknown symbolic knowledge is then the actual target to be learned. A classical example is *program synthesis*, where the goal is to learn the program from positive and negative examples of the desired input-output behaviour. Ideally, all positive pairs and none of the negatives should be covered. Many systems learn logic programs, i.e. NTPs

---

[3]Even though, many of them use a generative model to tackle the classification task, instead of a conditional one.

[102], $\partial$ILP [39], DeepProbLog[72], NeuralLP[137], DiffLog[112], DeepCoder[6]. Sometimes, the input-output pairs are not part of the training dataset, but are actually generated by a black-box neural model. The induced programs then explain the behaviour of the model, which relates NeSy to the domain of *explainability* [17].

## 9. Open Challenges

To conclude, we list some interesting challenges for NeSy.

*Semantics.* The statistical relational AI and probabilistic graphical model communities have devoted a lot of attention to the semantics of its models. This has resulted in several clear choices (such as directed vs. undirected, trace-based vs. possible world [104]), with corresponding strengths and weaknesses that clarify the relationships between the different models. Workshops have been held on this topic[4]. Furthermore, some researchers have investigated how to transform one type of model into another [59]. At the same time, weighted model counting has emerged as a common assembly language for inference. The situation in neurosymbolic computation today is very much like that of the early days in statistical relational learning, in which there were many competing formalisms, sometimes characterized as the statistical relational learning alphabet soup. It would be great to get more insight into the semantics of neurosymbolic approaches and their relationships. This survey hopes to contribute towards this goal.

*Probabilistic reasoning.* Although relatively few methods explore the integration of logical and neural methods from a probabilistic perspective, we believe that a probabilistic approach is a very natural way to integrate the two, since it has been shown [98] how one can recover the single methods as special cases. However, many open questions remain. Probabilistic inference is computationally expensive, usually requiring approximations. It would be interesting to determine exactly how probabilistic approximate inference compares with other approximations based on relaxations of the logic, like fuzzy logic.

*Fuzzy semantics.* The selection of the t-norm fuzzy logic and the corresponding translation of the connectives is very heterogeneous in the literature. It is often unclear which properties of Boolean logic a model is preserving, while there is a tendency to consider fuzzy logic as a continuous surrogate of Boolean logic without considering implications for the semantics. There is a clear need for further studies in this field. On the one hand, one may want to define new models which are natively fuzzy, thus not requiring a translation from Boolean logic. On the other hand, an interesting research direction concerns the characterisation of what are appropriate fuzzy approximations of Boolean logic relative to a set of properties that one wants to preserve (see Section 4).

---

[4]For instance, https://pps2018.luddy.indiana.edu/

*Structure learning.* While significant progress has been made on learning the structure of purely relational models (without probabilities), learning StarAI models remains a major challenge due to the complexity of inference and the combinatorial nature of the problem. Incorporating neural aspects complicates the problem even more. NeSy methods have certainly shown potential for addressing this problem (Section 5), but the existing methods are still limited and mostly domain-specific which impedes their wide application. For instance, the current systems that support structure learning require user effort to specify the clause templates or write a sketch of a model.

*Scaling inference.* Scalable inference is a major challenge for StarAI and therefore also for NeSy approaches with an explicit logical or probabilistic reasoning component. Investigating to what extent neural methods can help with this challenge by means of lifted (exploiting symmetries in models) or approximate inference, as well as reasoning from intermediate representations [1], are promising future research directions.

*Data efficiency.* A major advantage of StarAI methods, as compared to neural ones, is their data efficiency – StarAI methods can efficiently learn from small amounts of data, whereas neural methods are data hungry. On the other hand, StarAI methods do not scale to big data sets, while neural methods can easily handle them. We believe that understanding how these methods can help each other to overcome their complementary weaknesses, is a promising research direction.

*Symbolic representation learning.* The effectiveness of deep learning comes from the ability to change the representation of the data so that the target task becomes easier to solve. The ability to change the representation also at the symbolic level would significantly increase the capabilities of NeSy systems. This is a major open challenge for which neurally inspired methods could help achieve progress [19, 36].

**Acknowledgements**

# References

[1] Ralph Abboud, İsmail İlkan Ceylan, and Thomas Lukasiewicz. Learning to reason: Leveraging neural networks for approximate DNF counting. In *34th Conference on Artificial Intelligence, AAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 3097–3104. AAAI Press, 2020.

[2] Marjan Alirezaie, Martin Längkvist, Michael Sioutis, and Amy Loutfi. Semantic referee: A neural-symbolic framework for enhancing geospatial semantic segmentation. *Semantic Web*, 10(5):863–880, 2019.

[3] Stephen H. Bach, Matthias Broecheler, Bert Huang, and Lise Getoor. Hinge-loss markov random fields and probabilistic soft logic. *J. Mach. Learn. Res.*, 18:109:1–109:67, 2017.

[4] Sebastian Bader and Pascal Hitzler. Dimensions of neural-symbolic integration - A structured survey. In Sergei N. Artëmov, Howard Barringer, Artur S. d'Avila Garcez, Luís C. Lamb, and John Woods, editors, *We Will Show Them! Essays in Honour of Dov Gabbay, Volume One*, pages 167–194. College Publications, 2005.

[5] Samy Badreddine, Artur S. d'Avila Garcez, Luciano Serafini, and Michael Spranger. Logic tensor networks. *Artif. Intell.*, 303:103649, 2022.

[6] Matej Balog, Alexander L. Gaunt, Marc Brockschmidt, Sebastian Nowozin, and Daniel Tarlow. Deepcoder: Learning to write programs. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017*. OpenReview.net, 2017.

[7] Pablo Barceló, Egor V. Kostylev, Mikaël Monet, Jorge Pérez, Juan L. Reutter, and Juan Pablo Silva. The logical expressiveness of graph neural networks. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.

[8] Peter W. Battaglia, Jessica B. Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinícius Flores Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, Çaglar Gülçehre, H. Francis Song, Andrew J. Ballard, Justin Gilmer, George E. Dahl, Ashish Vaswani, Kelsey R. Allen, Charles Nash, Victoria Langston, Chris Dyer, Nicolas Heess, Daan Wierstra, Pushmeet Kohli, Matthew Botvinick, Oriol Vinyals, Yujia Li, and Razvan Pascanu. Relational inductive biases, deep learning, and graph networks. *CoRR*, abs/1806.01261, 2018.

[9] Radim Bělohlávek, Joseph W Dauben, and George J Klir. *Fuzzy logic and mathematics: a historical perspective*. Oxford University Press, 2017.

[10] Yoshua Bengio and Gary Marcus. Ai debate: The best way forward for ai https://montrealartificialintelligence.com/aidebate/, (checked on 13 december 2021), 2020.

[11] Tarek R. Besold, Artur d'Avila Garcez, and Luis C. Lamb. Human-Like Neural-Symbolic Computing (Dagstuhl Seminar 17192). *Dagstuhl Reports*, 7(5):56–83, 2017.

[12] Tarek R. Besold, Artur S. d'Avila Garcez, Sebastian Bader, Howard Bowman, Pedro M. Domingos, Pascal Hitzler, Kai-Uwe Kühnberger, Luís C. Lamb, Daniel Lowd, Priscila Machado Vieira Lima, Leo de Penning, Gadi Pinkas, Hoifung Poon, and Gerson Zaverucha. Neural-symbolic learning and reasoning: A survey and interpretation. *CoRR*, abs/1711.03902, 2017.

[13] Antoine Bordes, Nicolas Usunier, Alberto García-Durán, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. In *Advances in Neural Information Processing Systems 27, NeurIPS 2013, December 5-8, 2013, Lake Tahoe, Nevada*, pages 2787–2795, 2013.

[14] Matko Bosnjak, Tim Rocktäschel, Jason Naradowsky, and Sebastian Riedel. Programming with a differentiable forth interpreter. In Doina Precup and Yee Whye Teh, editors, *34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 547–556. PMLR, 2017.

[15] Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien. *Semi-supervised learning*. The MIT Press Cambridge, MA, USA:, 2006.

[16] Swarat Chaudhuri, Kevin Ellis, Oleksandr Polozov, Rishabh Singh, Armando Solar-Lezama, and Yisong Yue. Neurosymbolic programming. *Found. Trends Program. Lang.*, 7(3):158–243, 2021.

[17] Gabriele Ciravegna, Pietro Barbiero, Francesco Giannini, Marco Gori, Pietro Liò, Marco Maggini, and Stefano Melacci. Logic explained networks. *Artif. Intell.*, 314:103822, 2023.

[18] William W. Cohen, Fan Yang, and Kathryn Mazaitis. Tensorlog: Deep learning meets probabilistic dbs. *CoRR*, abs/1707.05390, 2017.

[19] Andrew Cropper. Playgol: Learning programs through play. In Sarit Kraus, editor, *28th International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019*, pages 6074–6080. ijcai.org, 2019.

[20] Andrew Cropper and Stephen H. Muggleton. Metagol system. https://github.com/metagol/metagol, 2016.

[21] James Cussens. Parameter estimation in stochastic logic programs. *Mach. Learn.*, 44(3):245–271, 2001.

[22] Wang-Zhou Dai, Qiu-Ling Xu, Yang Yu, and Zhi-Hua Zhou. Bridging machine learning and logical reasoning by abductive learning. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural*

*Information Processing Systems 32, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 2811–2822, 2019.

[23] Adnan Darwiche. On the tractable counting of theory models and its application to truth maintenance and belief revision. *J. Appl. Non Class. Logics*, 11(1-2):11–34, 2001.

[24] Adnan Darwiche. SDD: A new canonical representation of propositional knowledge bases. In Toby Walsh, editor, *22nd International Joint Conference on Artificial Intelligence, IJCAI 2011, Barcelona, Catalonia, Spain, July 16-22, 2011*, pages 819–826. IJCAI/AAAI, 2011.

[25] Adnan Darwiche and Pierre Marquis. A knowledge compilation map. *Journal of Artificial Intelligence Research*, 17:229–264, 2002.

[26] Tirtharaj Dash, Sharad Chitlangia, Aditya Ahuja, and Ashwin Srinivasan. How to tell deep neural networks what we know. *CoRR*, abs/2107.10295, 2021.

[27] Artur S. d'Avila Garcez, Marco Gori, Pascal Hitzler, and Luís C. Lamb. Neural-symbolic learning and reasoning (dagstuhl seminar 14381). *Dagstuhl Reports*, 4(9):50–84, 2014.

[28] Artur S. d'Avila Garcez, Marco Gori, Luís C. Lamb, Luciano Serafini, Michael Spranger, and Son N. Tran. Neural-symbolic computing: An effective methodology for principled integration of machine learning and reasoning. *FLAP*, 6(4):611–632, 2019.

[29] Artur S. d'Avila Garcez and Luís C. Lamb. A connectionist computational model for epistemic and temporal reasoning. *Neural Comput.*, 18(7):1711–1738, 2006.

[30] Artur S. d'Avila Garcez, Luís C. Lamb, and Dov M. Gabbay. Connectionist modal logic: Representing modalities in neural networks. *Theor. Comput. Sci.*, 371(1-2):34–53, 2007.

[31] Luc De Raedt, Angelika Kimmig, and Hannu Toivonen. Problog: A probabilistic prolog and its application in link discovery. In *20th International Joint Conference on Artificial Intelligence, IJCAI 2007, Hyderabad, India, January 6-12, 2007*, pages 2462–2467, 2007.

[32] Thomas Demeester, Tim Rocktäschel, and Sebastian Riedel. Lifted rule injection for relation embeddings. In Jian Su, Xavier Carreras, and Kevin Duh, editors, *Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, pages 1389–1399. The Association for Computational Linguistics, 2016.

[33] Michelangelo Diligenti, Marco Gori, and Claudio Saccà. Semantic-based regularization for learning and inference. *Artif. Intell.*, 244:143–165, 2017.

[34] Ivan Donadello, Luciano Serafini, and Artur S. d'Avila Garcez. Logic tensor networks for semantic image interpretation. In Carles Sierra, editor, *26th International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, pages 1596–1602. ijcai.org, 2017.

[35] Honghua Dong, Jiayuan Mao, Tian Lin, Chong Wang, Lihong Li, and Denny Zhou. Neural logic machines. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.

[36] Sebastijan Dumancic, Tias Guns, Wannes Meert, and Hendrik Blockeel. Learning relational representations with auto-encoding logic programs. In Sarit Kraus, editor, *28th International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019*, pages 6081–6087. ijcai.org, 2019.

[37] Kevin Ellis, Lucas Morales, Mathias Sablé-Meyer, Armando Solar-Lezama, and Josh Tenenbaum. Learning libraries of subroutines for neurally-guided bayesian program induction. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 7816–7826, 2018.

[38] Kevin Ellis, Maxwell I. Nye, Yewen Pu, Felix Sosa, Josh Tenenbaum, and Armando Solar-Lezama. Write, execute, assess: Program synthesis with a REPL. *CoRR*, abs/1906.04604, 2019.

[39] Richard Evans and Edward Grefenstette. Learning explanatory rules from noisy data. *J. Artif. Intell. Res.*, 61:1–64, 2018.

[40] Marc Fischer, Mislav Balunovic, Dana Drachsler-Cohen, Timon Gehr, Ce Zhang, and Martin T. Vechev. DL2: training and querying neural networks with logic. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 1931–1941. PMLR, 2019.

[41] Peter A. Flach. *Simply logical - intelligent reasoning by example*. Wiley professional computing. Wiley, 1994.

[42] Paolo Frasconi, Fabrizio Costa, Luc De Raedt, and Kurt De Grave. klog: A language for logical and relational learning with kernels. *Artif. Intell.*, 217:117–143, 2014.

[43] Nir Friedman, Lise Getoor, Daphne Koller, and Avi Pfeffer. Learning probabilistic relational models. In Thomas Dean, editor, *16th International Joint Conference on Artificial Intelligence, IJCAI 99, Stockholm, Sweden, July 31 - August 6, 1999. 2 Volumes, 1450 pages*, pages 1300–1309. Morgan Kaufmann, 1999.

[44] L. Getoor and B. Taskar, editors. *An Introduction to Statistical Relational Learning.* MIT Press, 2007.

[45] Francesco Giannini, Michelangelo Diligenti, Marco Gori, and Marco Maggini. On a convex logic fragment for learning and reasoning. *IEEE Trans. Fuzzy Syst.*, 27(7):1407–1416, 2019.

[46] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. Neural message passing for quantum chemistry. In Doina Precup and Yee Whye Teh, editors, *34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 1263–1272. PMLR, 2017.

[47] Martin Grohe. The logic of graph neural networks. In *36th Annual ACM/IEEE Symposium on Logic in Computer Science, LICS 2021, Rome, Italy, June 29 - July 2, 2021*, pages 1–17. IEEE, 2021.

[48] Shu Guo, Quan Wang, Lihong Wang, Bin Wang, and Li Guo. Jointly embedding knowledge graphs and logical rules. In Jian Su, Xavier Carreras, and Kevin Duh, editors, *Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, pages 192–202. The Association for Computational Linguistics, 2016.

[49] Bernd Gutmann, Angelika Kimmig, Kristian Kersting, and Luc De Raedt. Parameter learning in probabilistic databases: A least squares approach. In Walter Daelemans, Bart Goethals, and Katharina Morik, editors, *European Conference on Machine Learning and Knowledge Discovery in Databases, ECML/PKDD 2008, Antwerp, Belgium, September 15-19, 2008, Proceedings, Part I*, volume 5211 of *Lecture Notes in Computer Science*, pages 473–488. Springer, 2008.

[50] Bernd Gutmann, Ingo Thon, and Luc De Raedt. Learning the parameters of probabilistic logic programs from interpretations. In Dimitrios Gunopulos, Thomas Hofmann, Donato Malerba, and Michalis Vazirgiannis, editors, *European Conference on Machine Learning, ECML PKDD 2011, Athens, Greece, September 5-9, 2011. Proceedings, Part I*, volume 6911 of *Lecture Notes in Computer Science*, pages 581–596. Springer, 2011.

[51] Christopher Hahn, Frederik Schmitt, Jens U. Kreber, Markus Norman Rabe, and Bernd Finkbeiner. Teaching temporal logics to neural networks. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.

[52] Joseph Y. Halpern. An analysis of first-order logics of probability. *Artif. Intell.*, 46(3):311–350, 1990.

[53] Joseph Y. Halpern. *Reasoning about uncertainty.* MIT Press, 2005.

[54] William L. Hamilton, Payal Bajaj, Marinka Zitnik, Dan Jurafsky, and Jure Leskovec. Embedding logical queries on knowledge graphs. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 2030–2041, 2018.

[55] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. *CoRR*, abs/1503.02531, 2015.

[56] Sepp Hochreiter. Toward a broad AI. *Commun. ACM*, 65(4):56–57, 2022.

[57] Zhiting Hu, Xuezhe Ma, Zhengzhong Liu, Eduard H. Hovy, and Eric P. Xing. Harnessing deep neural networks with logic rules. In *54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*. The Association for Computer Linguistics, 2016.

[58] Jiani Huang, Ziyang Li, Binghong Chen, Karan Samel, Mayur Naik, Le Song, and Xujie Si. Scallop: From probabilistic deductive databases to scalable differentiable reasoning. pages 25134–25145, 2021.

[59] Manfred Jaeger. Model-theoretic expressivity analysis. In Luc De Raedt, Paolo Frasconi, Kristian Kersting, and Stephen H. Muggleton, editors, *Probabilistic Inductive Logic Programming - Theory and Applications*, volume 4911 of *Lecture Notes in Computer Science*, pages 325–339. Springer, 2008.

[60] Ashwin Kalyan, Abhishek Mohta, Oleksandr Polozov, Dhruv Batra, Prateek Jain, and Sumit Gulwani. Neural-guided deductive search for real-time program synthesis from examples. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018*. OpenReview.net, 2018.

[61] Kristian Kersting and Luc De Raedt. Bayesian logic programming: Theory and tool. In L. Getoor and B. Taskar, editors, *An introduction to Statistical Relational Learning*. MIT Press, 2007.

[62] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. 2015.

[63] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017*. OpenReview.net, 2017.

[64] Stanley Kok and Pedro M. Domingos. Learning the structure of markov logic networks. In Luc De Raedt and Stefan Wrobel, editors, *22nd International Conference on Machine Learning, (ICML 2005), Bonn, Germany,*

*August 7-11, 2005*, volume 119 of *ACM International Conference Proceeding Series*, pages 441–448. ACM, 2005.

[65] Stanley Kok and Pedro M. Domingos. Learning markov logic networks using structural motifs. In Johannes Fürnkranz and Thorsten Joachims, editors, *27th International Conference on Machine Learning (ICML-10), June 21-24, 2010, Haifa, Israel*, pages 551–558. Omnipress, 2010.

[66] Daphne Koller and Nir Friedman. *Probabilistic Graphical Models - Principles and Techniques*. MIT Press, 2009.

[67] Luís C. Lamb, Artur S. d'Avila Garcez, Marco Gori, Marcelo O. R. Prates, Pedro H. C. Avelar, and Moshe Y. Vardi. Graph neural networks meet neural-symbolic computing: A survey and perspective. In Christian Bessiere, editor, *29th International Joint Conference on Artificial Intelligence, IJCAI 2020*, pages 4877–4884. ijcai.org, 2020.

[68] Tao Li and Vivek Srikumar. Augmenting neural networks with first-order logic. *CoRR*, abs/1906.06298, 2019.

[69] John W. Lloyd. *Foundations of Logic Programming, 2nd Edition*. Springer, 1987.

[70] Daniel Lowd and Pedro M. Domingos. Efficient weight learning for markov logic networks. In Joost N. Kok, Jacek Koronacki, Ramón López de Mántaras, Stan Matwin, Dunja Mladenic, and Andrzej Skowron, editors, *11th European Conference on Principles and Practice of Knowledge Discovery in Databases, PKDD 2007, Warsaw, Poland, September 17-21, 2007, Proceedings*, volume 4702 of *Lecture Notes in Computer Science*, pages 200–211. Springer, 2007.

[71] Jayanta Mandi, Víctor Bucarey, Maxime Mulamba Ke Tchomba, and Tias Guns. Decision-focused learning: Through the lens of learning to rank. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato, editors, *39th International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 14935–14947. PMLR, 2022.

[72] Robin Manhaeve, Sebastijan Dumancic, Angelika Kimmig, Thomas Demeester, and Luc De Raedt. Deepproblog: Neural probabilistic logic programming. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 3753–3763, 2018.

[73] Robin Manhaeve, Giuseppe Marra, and Luc De Raedt. Approximate inference for neural probabilistic logic programming. In Meghyn Bienvenu,

Gerhard Lakemeyer, and Esra Erdem, editors, *18th International Conference on Principles of Knowledge Representation and Reasoning, KR 2021, Online event, November 3-12, 2021*, pages 475–486, 2021.

[74] Jiayuan Mao, Chuang Gan, Pushmeet Kohli, Joshua B. Tenenbaum, and Jiajun Wu. The neuro-symbolic concept learner: Interpreting scenes, words, and sentences from natural supervision. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.

[75] Giuseppe Marra. *Bridging symbolic and subsymbolic reasoning with Mini-Max Entropy models.* PhD thesis, University of Florence, 2 2020.

[76] Giuseppe Marra, Michelangelo Diligenti, Francesco Giannini, Marco Gori, and Marco Maggini. Relational neural machines. In Giuseppe De Giacomo, Alejandro Catalá, Bistra Dilkina, Michela Milano, Senén Barro, Alberto Bugarín, and Jérôme Lang, editors, *24th European Conference on Artificial Intelligence, ECAI 2020, 29 August-8 September 2020, Santiago de Compostela, Spain, August 29 - September 8, 2020*, volume 325 of *Frontiers in Artificial Intelligence and Applications*, pages 1340–1347. IOS Press, 2020.

[77] Giuseppe Marra, Francesco Giannini, Michelangelo Diligenti, and Marco Gori. Integrating learning and reasoning with deep logic models. In Ulf Brefeld, Élisa Fromont, Andreas Hotho, Arno J. Knobbe, Marloes H. Maathuis, and Céline Robardet, editors, *European Conference on Machine Learning and Knowledge Discovery in Databases, ECML PKDD 2019, Würzburg, Germany, September 16-20, 2019, Proceedings, Part II*, volume 11907 of *Lecture Notes in Computer Science*, pages 517–532. Springer, 2019.

[78] Giuseppe Marra and Ondrej Kuzelka. Neural markov logic networks. In Cassio P. de Campos, Marloes H. Maathuis, and Erik Quaeghebeur, editors, *37th Conference on Uncertainty in Artificial Intelligence, UAI 2021, Virtual Event, 27-30 July 2021*, volume 161 of *Proceedings of Machine Learning Research*, pages 908–917. AUAI Press, 2021.

[79] Pasquale Minervini, Matko Bosnjak, Tim Rocktäschel, Sebastian Riedel, and Edward Grefenstette. Differentiable reasoning on large knowledge bases and natural language. In *34th Conference on Artificial Intelligence, AAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 5182–5190. AAAI Press, 2020.

[80] Pasquale Minervini, Thomas Demeester, Tim Rocktäschel, and Sebastian Riedel. Adversarial sets for regularising neural link predictors. In Gal Elidan, Kristian Kersting, and Alexander T. Ihler, editors, *33rd Conference on Uncertainty in Artificial Intelligence, UAI 2017, Sydney, Australia, August 11-15, 2017*. AUAI Press, 2017.

[81] Pasquale Minervini, Sebastian Riedel, Pontus Stenetorp, Edward Grefen-stette, and Tim Rocktäschel. Learning reasoning strategies in end-to-end differentiable proving. In *37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 6938–6949. PMLR, 2020.

[82] Eleonora Misino, Giuseppe Marra, and Emanuele Sansone. VAEL: bridging variational autoencoders and probabilistic logic programming. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35, NeurIPS 2022, November 29 December 4, 2022, New Orleans, Luisiana*, pages 4667–4679, 2022.

[83] Tom M. Mitchell. Generalization as search. *Artif. Intell.*, 18(2):203–226, 1982.

[84] Christopher Morris, Martin Ritzert, Matthias Fey, William L. Hamilton, Jan Eric Lenssen, Gaurav Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural networks. In *33rd Conference on Artificial Intelligence, AAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pages 4602–4609. AAAI Press, 2019.

[85] Stephen Muggleton. Stochastic logic programs. *Advances in inductive logic programming*, 32, 1996.

[86] Stephen Muggleton and Luc De Raedt. Inductive logic programming: Theory and methods. *J. Log. Program.*, 19/20:629–679, 1994.

[87] Deepak Nathani, Jatin Chauhan, Charu Sharma, and Manohar Kaul. Learning attention-based embeddings for relation prediction in knowledge graphs. In Anna Korhonen, David R. Traum, and Lluís Màrquez, editors, *57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 4710–4723. Association for Computational Linguistics, 2019.

[88] Maximilian Nickel, Kevin Murphy, Volker Tresp, and Evgeniy Gabrilovich. A review of relational machine learning for knowledge graphs. *Proc. IEEE*, 104(1):11–33, 2016.

[89] Maximilian Nickel, Volker Tresp, and Hans-Peter Kriegel. A three-way model for collective learning on multi-relational data. In Lise Getoor and Tobias Scheffer, editors, *28th International Conference on Machine Learning, ICML 2011, Bellevue, Washington, USA, June 28 - July 2, 2011*, pages 809–816. Omnipress, 2011.

[90] Maxwell I. Nye, Armando Solar-Lezama, Josh Tenenbaum, and Brenden M. Lake. Learning compositional rules via neural program synthesis. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.

[91] Judea Pearl. *Probabilistic reasoning in intelligent systems - networks of plausible inference*. Morgan Kaufmann series in representation and reasoning. Morgan Kaufmann, 1989.

[92] David Poole. Probabilistic horn abduction and bayesian networks. *Artif. Intell.*, 64(1):81–129, 1993.

[93] Meng Qu and Jian Tang. Probabilistic logic neural networks for reasoning. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 7710–7720, 2019.

[94] Luc De Raedt. *Logical and relational learning*. Cognitive Technologies. Springer, 2008.

[95] Luc De Raedt, Anton Dries, Ingo Thon, Guy Van den Broeck, and Mathias Verbeke. Inducing probabilistic relational rules from probabilistic examples. In Qiang Yang and Michael J. Wooldridge, editors, *24th International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina, July 25-31, 2015*, pages 1835–1843. AAAI Press, 2015.

[96] Luc De Raedt, Kristian Kersting, Sriraam Natarajan, and David Poole. *Statistical Relational Artificial Intelligence: Logic, Probability, and Computation*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2016.

[97] Luc De Raedt and Angelika Kimmig. Probabilistic (logic) programming concepts. *Mach. Learn.*, 100(1):5–47, 2015.

[98] Luc De Raedt, Robin Manhaeve, Sebastijan Dumancic, Thomas Demeester, and Angelika Kimmig. Neuro-symbolic = neural + logical + probabilistic. In Derek Doran, Artur S. d'Avila Garcez, and Freddy Lécué, editors, *International Workshop on Neural-Symbolic Learning and Reasoning (NeSy 2019), Annual workshop of the Neural-Symbolic Learning and Reasoning Association, Macao, China, August 12, 2019*, 2019.

[99] Hongyu Ren and Jure Leskovec. Beta embeddings for multi-hop logical reasoning in knowledge graphs. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.

[100] Matthew Richardson and Pedro M. Domingos. Markov logic networks. *Mach. Learn.*, 62(1-2):107–136, 2006.

[101] Ryan Riegel, Alexander G. Gray, Francois P. S. Luus, Naweed Khan, Ndivhuwo Makondo, Ismail Yunus Akhalwaya, Haifeng Qian, Ronald Fagin, Francisco Barahona, Udit Sharma, Shajith Ikbal, Hima Karanam,

Sumit Neelam, Ankita Likhyani, and Santosh K. Srivastava. Logical neural networks. *CoRR*, abs/2006.13155, 2020.

[102] Tim Rocktäschel and Sebastian Riedel. End-to-end differentiable proving. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 30, NeurIPA 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 3788–3800, 2017.

[103] Tim Rocktäschel, Sameer Singh, and Sebastian Riedel. Injecting logical background knowledge into embeddings for relation extraction. In Rada Mihalcea, Joyce Yue Chai, and Anoop Sarkar, editors, *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL HLT 2015 , Denver, Colorado, USA, May 31 - June 5, 2015*, pages 1119–1129. The Association for Computational Linguistics, 2015.

[104] Stuart J. Russell. Unifying logic and probability. *Commun. ACM*, 58(7):88–97, 2015.

[105] Taisuke Sato. A statistical learning method for logic programs with distribution semantics. In Leon Sterling, editor, *12nd International Conference on Logic Programming, Tokyo, Japan, June 13-16, 1995*, pages 715–729. MIT Press, 1995.

[106] Taisuke Sato and Yoshitaka Kameya. PRISM: A language for symbolic-statistical modeling. In *15th International Joint Conference on Artificial Intelligence, IJCAI 97, Nagoya, Japan, August 23-29, 1997, 2 Volumes*, pages 1330–1339. Morgan Kaufmann, 1997.

[107] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE transactions on neural networks*, 20(1):61–80, 2008.

[108] Michael Sejr Schlichtkrull, Thomas N. Kipf, Peter Bloem, Rianne van den Berg, Ivan Titov, and Max Welling. Modeling relational data with graph convolutional networks. In Aldo Gangemi, Roberto Navigli, Maria-Esther Vidal, Pascal Hitzler, Raphaël Troncy, Laura Hollink, Anna Tordai, and Mehwish Alam, editors, *15th International Conference on the The Semantic Web, ESWC 2018, Heraklion, Crete, Greece, June 3-7, 2018, Proceedings*, volume 10843 of *Lecture Notes in Computer Science*, pages 593–607. Springer, 2018.

[109] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Gallagher, and Tina Eliassi-Rad. Collective classification in network data. *AI Mag.*, 29(3):93–106, 2008.

[110] Chao Shang, Yun Tang, Jing Huang, Jinbo Bi, Xiaodong He, and Bowen Zhou. End-to-end structure-aware convolutional networks for knowledge

base completion. In *33rd Conference on Artificial Intelligence, AAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pages 3060–3067. AAAI Press, 2019.

[111] Hikaru Shindo, Viktor Pfanschilling, Devendra Singh Dhami, and Kristian Kersting. $\alpha$ILP: thinking visual scenes as differentiable logic programs. *Mach. Learn.*, 112(5):1465–1497, 2023.

[112] Xujie Si, Mukund Raghothaman, Kihong Heo, and Mayur Naik. Synthesizing datalog programs using numerical relaxation. In Sarit Kraus, editor, *28th International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019*, pages 6117–6124. ijcai.org, 2019.

[113] Arseny Skryagin, Wolfgang Stammer, Daniel Ochs, Devendra Singh Dhami, and Kristian Kersting. Neural-probabilistic answer set programming. In Gabriele Kern-Isberner, Gerhard Lakemeyer, and Thomas Meyer, editors, *19th International Conference on Principles of Knowledge Representation and Reasoning, KR 2022, Haifa, Israel, July 31 - August 5, 2022*, 2022.

[114] Arseny Skryagin, Karl Stelzner, Alejandro Molina, Fabrizio Ventola, Zhongjie Yu, and Kristian Kersting. Sum-product logic: integrating probabilistic circuits into deepproblog. 2020.

[115] Lennert De Smet, Pedro Zuidberg Dos Martires, Robin Manhaeve, Giuseppe Marra, Angelika Kimmig, and Luc De Raedt. Neural probabilistic logic programming in discrete-continuous domains. In Robin J. Evans and Ilya Shpitser, editors, *Uncertainty in Artificial Intelligence, UAI 2023, July 31 - 4 August 2023, Pittsburgh, PA, USA*, volume 216 of *Proceedings of Machine Learning Research*, pages 529–538. PMLR, 2023.

[116] Gustav Sourek, Vojtech Aschenbrenner, Filip Zelezný, Steven Schockaert, and Ondrej Kuzelka. Lifted relational neural networks: Efficient learning of latent relational structures. *J. Artif. Intell. Res.*, 62:69–100, 2018.

[117] Gustav Sourek, Filip Zelezný, and Ondrej Kuzelka. Beyond graph neural networks with lifted relational neural networks. *Mach. Learn.*, 110(7):1695–1738, 2021.

[118] Naoya Takeishi and Kosuke Akimoto. Knowledge-based distant regularization in learning probabilistic models. *CoRR*, abs/1806.11332, 2018.

[119] Geoffrey G. Towell and Jude W. Shavlik. Knowledge-based artificial neural networks. *Artif. Intell.*, 70(1-2):119–165, 1994.

[120] Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. Complex embeddings for simple link prediction. In Maria-Florina Balcan and Kilian Q. Weinberger, editors, *33nd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 2071–2080. JMLR.org, 2016.

[121] Efthymia Tsamoura, Timothy M. Hospedales, and Loizos Michael. Neural-symbolic integration: A compositional perspective. In *35th Conference on Artificial Intelligence, AAAI 2021, Virtual Event, February 2-9, 2021*, pages 5051–5060. AAAI Press, 2021.

[122] Lazar Valkov, Dipak Chaudhari, Akash Srivastava, Charles Sutton, and Swarat Chaudhuri. HOUDINI: lifelong learning as program synthesis. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 8701–8712, 2018.

[123] Michael van Bekkum, Maaike de Boer, Frank van Harmelen, André Meyer-Vitali, and Annette ten Teije. Modular design patterns for hybrid learning and reasoning systems. *Appl. Intell.*, 51(9):6528–6546, 2021.

[124] Emile van Krieken, Erman Acar, and Frank van Harmelen. Analyzing differentiable fuzzy logic operators. *Artif. Intell.*, 302:103602, 2022.

[125] Moshe Y. Vardi. Why is modal logic so robustly decidable? In Neil Immerman and Phokion G. Kolaitis, editors, *Descriptive Complexity and Finite Models, DIMACS Workshop 1996, Princeton, New Jersey, USA, January 14-17, 1996*, volume 31 of *DIMACS Series in Discrete Mathematics and Theoretical Computer Science*, pages 149–183. DIMACS/AMS, 1996.

[126] Shikhar Vashishth, Soumya Sanyal, Vikram Nitin, and Partha P. Talukdar. Composition-based multi-relational graph convolutional networks. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.

[127] Po-Wei Wang, Priya L. Donti, Bryan Wilder, and J. Zico Kolter. Sat-net: Bridging deep learning and logical reasoning using a differentiable satisfiability solver. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 6545–6554. PMLR, 2019.

[128] Quan Wang, Zhendong Mao, Bin Wang, and Li Guo. Knowledge graph embedding: A survey of approaches and applications. *IEEE Trans. Knowl. Data Eng.*, 29(12):2724–2743, 2017.

[129] Wenya Wang and Sinno Jialin Pan. Integrating deep learning with logic fusion for information extraction. pages 9225–9232, 2020.

[130] Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. Knowledge graph embedding by translating on hyperplanes. In Carla E. Brodley and Peter Stone, editors, *28th Conference on Artificial Intelligence, AAAI 2014, July 27 -31, 2014, Québec City, Québec, Canada*, pages 1112–1119. AAAI Press, 2014.

[131] Leon Weber, Pasquale Minervini, Jannes Münchmeyer, Ulf Leser, and Tim Rocktäschel. Nlprolog: Reasoning with weak unification for question answering in natural language. In Anna Korhonen, David R. Traum, and Lluís Màrquez, editors, *57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 6151–6161. Association for Computational Linguistics, 2019.

[132] Thomas Winters, Giuseppe Marra, Robin Manhaeve, and Luc De Raedt. Deepstochlog: Neural stochastic logic programming. In *36th Conference on Artificial Intelligence, AAAI 2022, Virtual Event, February 22 - March 1, 2022*, pages 10090–10100. AAAI Press, 2022.

[133] Jingyi Xu, Zilu Zhang, Tal Friedman, Yitao Liang, and Guy Van den Broeck. A semantic loss function for deep learning with symbolic knowledge. In Jennifer G. Dy and Andreas Krause, editors, *35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 5498–5507. PMLR, 2018.

[134] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.

[135] Xiaoran Xu, Wei Feng, Yunsheng Jiang, Xiaohui Xie, Zhiqing Sun, and Zhi-Hong Deng. Dynamically pruned message passing networks for large-scale knowledge graph reasoning. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.

[136] Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. Embedding entities and relations for learning and inference in knowledge bases. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015*, 2015.

[137] Fan Yang, Zhilin Yang, and William W. Cohen. Differentiable learning of logical rules for knowledge base reasoning. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 30, NeurIPS 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 2319–2328, 2017.

[138] Zhun Yang, Adam Ishay, and Joohyung Lee. Neurasp: Embracing neural networks into answer set programming. In Christian Bessiere, editor, *29th International Joint Conference on Artificial Intelligence, IJCAI 2020*, pages 1755–1762. ijcai.org, 2020.

[139] Lotfi A Zadeh. Fuzzy logic and approximate reasoning. *Synthese*, 30(3-4):407–428, 1975.

[140] Lisa Zhang, Gregory Rosenblatt, Ethan Fetaya, Renjie Liao, William E. Byrd, Matthew Might, Raquel Urtasun, and Richard S. Zemel. Neural guided constraint logic programming for program synthesis. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 1744–1753, 2018.

[141] Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 5171–5181, 2018.

[142] Yuyu Zhang, Xinshi Chen, Yuan Yang, Arun Ramamurthy, Bo Li, Yuan Qi, and Le Song. Efficient probabilistic logic reasoning with graph neural networks. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.

[143] Pedro Zuidberg Dos Martires, Vincent Derkinderen, Robin Manhaeve, Wannes Meert, Angelika Kimmig, and Luc De Raedt. Transforming probabilistic programs into algebraic circuits for inference and learning. In *Program Transformations for ML Workshop at NeurIPS*, 2019.

## Appendix A. Knowledge graphs embeddings and Graph Neural Networks for neurosymbolic AI

In this appendix, we introduce two approaches, namely Knowledge Graph Embeddings and Graph Neural Networks, which are commonly used in relational tasks in the deep learning community. We analyze them in the spirit of the seven dimensions introduced in the paper. In fact, it turns out that they share many features with neurosymbolic systems. In this way, we would like to suggest that NeSy can also be found at the intersection of statistical relational and geometric deep learning approaches.

One of the most popular relational representations is that of a knowledge graph (KG). A KG is a multi-relational graph composed of entities (i.e. nodes) and relations (i.e. edges). It is common to represent an edge as a triple of the form: *(head entity, relation, tail entity)*, e.g. $(homer, fatherOf, bart)$.

StarAI has been extensively used [88] to solve many tasks on KGs: prediction of missing relationships (i.e. knowledge graph completion), prediction of properties of entities, or clustering entities based on their connectivity patterns. StarAI is particularly well suited to reasoning with knowledge graphs since it models explicitly the probabilistic dependencies among different relationships.

However, in order to scale to larger knowledge graphs, the probabilistic dependencies of StarAI models have been relaxed to give rise to a new class of scalable models based on latent features, which are particularly interesting from a neurosymbolic viewpoint. The key intuition behind relational latent feature models is that the relationships between entities can be more efficiently predicted by modeling simpler interactions in a latent feature space. Knowledge Graph Embeddings and Graph Neural Networks represent two of the ways to encode such latent representations.

### Appendix A.1. Knowledge Graph Embeddings

Knowledge Graph Embedding (KGE) models [128] assume that triples (i.e. relations) are conditionally independent given a set of global latent variables, called embeddings, for entities and relations. Therefore, the existence of a triple can be predicted by a scoring function $f(e_h, e_r, e_t)$, where $e_i$ is the embedding of the corresponding object.

KGE models mainly differ in terms of the scoring function and the embedding space they use. Translation or distance-based models [13, 130] use a scoring function measuring to what extent the tail of a triple can be obtained by a relation-specific translation of the head, i.e. $f(e_h, e_r, e_t) = ||e_h + e_r - e_t||$. Semantic matching methods [89] instead exploit similarity-based scoring functions, such as $f(e_h, W_r, e_t) = ||e_h W_r e_t^\top||$. It is interesting to note that standard multi-layer perceptrons are often used to learn this similarity measure, i.e. $f(e_h, W_r, e_t) = nn([e_h, e_t]; W_r)$. This is similar to neural interfaces typical of $lPN$ and $LpN$ models of Section 7, cf. the neural predicates of DeepProbLog in Example 13.

KGE learn from ground relations only. They learn the embeddings of entities and relations as to maximize the score for a set of known *True* triples. However, some methods incorporate higher level information, like first-order logical clauses [48, 32] or logical queries [54, 99], bridging KGE and multiple NeSy systems (like SL[133] or SBR[33]).

It is interesting to analyze KGE methods in terms of the dimensions we described in our paper. KGE methods mostly work as model-based, undirected methods. They constrain the embeddings to be coherent with the logical facts and rules, which are no longer used after learning. They are heavily based on subsymbolic representations. These models learn the correct parameters (i.e. embeddings) for the task at hand. Even if no explicit semantics needs to be given to the scoring function $f(e_h, W_r, e_t)$, a fuzzy logic interpretation is often used when injecting logical rules [48, 32]

### Appendix A.2. Graph Neural Networks

Graph Neural Networks (GNN) ([107, 134]) are deep learning models for dealing with graphs as input. Inference in these models can be cast as *message passing* [46, 8]: at each inference step (i.e. GNN layer), a node sends a message to all its outgoing neighbors and updates its state by aggregating all the messages coming from its ingoing neighbors. Messages are computed by standard neural networks. Each inference step computes a transformation of the representations

of the nodes. In the last layer, these representations are used to classify the nodes or are aggregated to classify the entire graph.

There are many connections between GNNs and neurosymbolic computation [67] since both of them apply neural methods to relational data: graphs for GNNs, and logic representations for NeSy. While many works try to close the gap between these two representations [42], the application of GNN techniques to logic-based graph structures is still very limited. The most related line of work is about using GNNs on knowledge graphs [141, 108, 110, 135, 87, 126]. The underlying idea is to differentiate the messages exchanged by two nodes if they are related by different kinds of edges. For example, given two nodes `homer` and `bart`, the message corresponding to the edge `fatherOf(homer,bart)` will be different from the one corresponding to `sameFamily(homer,bart)`. However, these models were intended as classifiers of known relational structures (nodes and edges) and not to reason about the knowledge graph itself (e.g. determining whether an edge exists between two nodes). To perform relational reasoning, GNN-based models rely on techniques from the KGE community on top of the representations extracted by the GNN. This often takes the shape of an auto-encoding scheme: a GNN encodes an input graph in a latent representation and a KGE-based factorization technique is used to reconstruct the whole graph [108].

An important characteristic of GNNs is that they rely exclusively on neural computation to perform inference (i.e. to compute messages) and there is no clear direction on how to inject external knowledge about inference, e.g. as logical rules. This contrasts with NeSy, where this is one of the main goals.

There are also some interesting connections between GNNs and StarAI models. In [93, 142], GNNs based on knowledge-graphs are not used as a modeling choice but rather to approximate inference in Markov Logic Networks, which is somewhat similar to regularization based methods (see Section 2). Similarly, in [1] GNNs are used to encode logical formulae expressed as graphs to approximate a weighted model counting problem.

Finally, it is interesting to analyse GNNs in the spirit of some of the dimensions of NeSy. GNNs act as directed models with a proof-based inference scheme: they perform a series of inference steps to compute the final answer. In the original version of GNNs [107], the node states are updated until a fixed point is reached, which resembles forward-chaining in logic programming. The representation of nodes belongs to a subsymbolic numerical space. Finally, GNNs can be considered as implicit structure learners: inference rules are learned through the learning of the neural message passing functions.

Graph Neural Networks have recently received a lot of attention from many different communities, thanks to the representation power of neural networks and the capability of learning in complex relational settings. It is no surprise that people have started to study the expressivity of this class of models. One of the most interesting analyses from a neurosymbolic viewpoint is measuring the expressivity of GNNs in terms of variable counting logics. Recently, [84, 7, 47] showed that GNNs are as expressive as 2-variable counting logic $C_2$. This fragment of first order logic admits formulas with at most two variables extended

with counting quantifiers. The expressivity of this fragment is limited compared to many neurosymbolic models, especially those based on logic. However, GNNs learn the logical structure of the problem implicitly as part of the message passing learning scheme and they rely neither on expert-provided knowledge nor on heavy combinatorial search strategies to structure learning (see Section 5). An open and challenging question that unites the GNN and NeSy communities is how to bring the expressivity to higher-order fragments [84], like in NeSy and StarAI, while keeping both the learning and the inference tractable, like in GNNs.

## Appendix  B. Fuzzy logic, fuzzyfication and soft-Satisfability

Fuzzy logic, as many-valued extension over Boolean logic, has a very long tradition [139, 9]. However, the use of fuzzy logic in StarAI and NeSy is not dictated by the need of dealing with vagueness, but by the advantageous computational properties of t-norms. Indeed, a common use case is to have an initial theory defined in Boolean logic which is *fuzzyfied*. Inference is then carried out with the fuzzyfied theory and the answers are eventually discretized back to Boolean values (usually using a threshold at 0.5).

The reason for this approach is that one would like to exploit the differentiability of t-norms to address logical inference of FOL theories in a more scalable way than standard combinatorial optimization algorithms (e.g. SAT solvers). This is particularly important in undirected and regularization-based methods (such as PSL [3] and LTN [5]). In fact, it has been shown [45] that there are fragments of fuzzy logic that can even provide convex inference problems. Example 23, however, shows that naively approaching logical inference through a fuzzy relaxation and gradient-based optimization can introduce unexpected behaviours.

---

EXAMPLE 23: FUZZYFICATION AND SOFT-SATISFABILITY

Let us consider a disjunction, like $A \vee B \vee C$. In Boolean logic, if we state that the disjunction is *satisfied* (i.e. *True*), then we expect at least one among the three variables to be *True*.

Suppose we want to find a truth assignment for all the variables that satisfies the disjunction above. The approach of the majority of NeSy fuzzy approaches is the following. First, the rule is relaxed into a fuzzy real function. For example, using the Łukasiewicz t-norm, $F_\oplus(A, B, C) = min(1, A + B + C)$. Secondly, a gradient-based algorithm (e.g. backpropagation with Adam [62]) is used to maximize the value of the formula with respect to the *fuzzy* truth degree of the three variables. Finally, the obtained *fuzzy* solution $A^\star, B^\star, C^\star$ is translated back into a Boolean assignment using a 0.5 threshold.

Let us consider a possible optimal fuzzy solution, like $(A^\star, B^\star, C^\star) = (0.34, 0.34, 0.34)$ and its discretized version $(A^\star, B^\star, C^\star) = (False, False, False)$, using a threshold at 0.5. The discretized solution does not satisfy the initial Boolean formula, even though it is a global optimum in the fuzzyfied

problem.

Similarly, [124] shows that, while it is very common to reason about universally quantified formulae in the form of $\forall x : A(x) \rightarrow B(x)$, like 'all humans are mortal', using gradients and fuzzy logic to make inference can be extremely counterintuitive, especially with specific t-norms such as the product t-norm. It is unclear whether there exists a generally accepted subset of properties of Boolean logic that one wants to preserve and whether one can define a t-norm that guarantees such properties.

*Appendix B.1. Distribution semantics and fuzzy logic semantics*

Another common reason for using fuzzy logic is to exploit a differentiable semantics. Then, gradient-based methods can be used to train the parameters of a weighted logical theory such as in LRNNs [116]. This contrasts with gradient-based training of the parameters of probabilistic logics based on the distribution semantics. Possible worlds in probabilistic logic are defined as possible assignments of truth values to all the ground atoms of a logical theory. The assignments of truth values specify the *semantics* of the logic. On the contrary, fuzzy logic assigns continuous truth degrees to formulas or proofs, which are *syntactic* structures. As a consequence, while the probability of an atom will always be equal to the sum of the probabilities of the worlds in which it is *True* (cf. Equation 1), the fuzzy degree of an atom may vary depending on how that atom has been proven or defined, as shown in Example 24.

EXAMPLE 24: DISTRIBUTION SEMANTICS VS FUZZY LOGIC SEMANTICS

Consider the following annotated program (adapted from [21]).

```
0.3::b.
a1 <- b.
a2 <- b,b.
```

Here, 0.3 is the label of b without any particular semantics yet. Given the idempotency of the Boolean conjunction, we would expect the scores of both a1 and a2 to be identical since both b and b ∧ b are *True* when b is *True*. In a probabilistic approach, the score is interpreted as a probability, i.e. $p(\mathtt{b}) = 0.3$. The probability of any atom is the sum of the probabilities of all the worlds where that atom is *True*. It is easy to see that, for both a1 and a2, these are the worlds where b is *True*, thus $p(\mathtt{a1}) = p(\mathtt{a2}) = p(\mathtt{b}) = 0.3$. On the other hand, in the fuzzy setting, the score of b is interpreted as its truth degree, i.e. $t(\mathtt{b}) = 0.3$. Let us consider the product t-norm, where $t(x \wedge y) = t(x)t(y)$, then $t(\mathtt{a1}) = t(\mathtt{b}) = 0.3$ while $t(\mathtt{a2}) = t(\mathtt{b})t(\mathtt{b}) = 0.09$. While this issue could be solved by choosing a different t-norm (e.g. the minimum t-norm), similar issues arise in different definitions.

The differences between the two semantics are not due to the probabilistic or the fuzzy semantics, but more to the distinction between semantics based

on possible worlds and semantics based on proofs or derivations. In fact, a similar behaviour is observed in Stochastic Logic Programs [21] under the name of *memoization.*