
Scallop: From Probabilistic Deductive Databases to Scalable Differentiable Reasoning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Deep learning and symbolic reasoning are complementary techniques for an intelligent
2 system. However, principled combinations of these techniques are typically
3 limited in scalability, rendering them ill-suited for real-world applications. We propose
4 Scallop, a system that builds upon probabilistic deductive databases, to bridge
5 this gap. The key insight underlying Scallop is a provenance framework that in-
6 troduces a tunable parameter to specify the level of reasoning granularity. Scallop
7 thereby i) generalizes exact probabilistic reasoning, ii) asymptotically reduces
8 computational cost, and iii) provides relative accuracy guarantees. On synthetic
9 tasks involving mathematical and logical reasoning, Scallop scales significantly
10 better without sacrificing accuracy compared to DeepProbLog, a principled neural
11 logic programming approach. Scallop also scales to a newly created real-world
12 Visual Question Answering (VQA) benchmark that requires multi-hop reasoning,
13 achieving 84.22% accuracy and outperforming two VQA-tailored models based on
14 Neural Module Networks and transformers by 12.42% and 21.66% respectively.

15 1 Introduction

16 Integrating deep learning and symbolic reasoning in a principled manner into a single effective system
17 is a fundamental problem in artificial intelligence [11]. Despite great potential in terms of accuracy,
18 interpretability, and generalizability, it is challenging to scale differentiable reasoning in the combined
19 system while preserving the benefits of the neural and symbolic sub-systems [27].

20 In this paper, we propose Scallop, a systematic and effective framework to address this problem.
21 The key insight underlying Scallop is a principled relaxation of exact probabilistic reasoning via a
22 parameter k that specifies the level of reasoning granularity. We observe that scalability is primarily
23 hindered by reasoning about *all* proofs in computing the probability of each outcome. For a given k ,
24 Scallop only reasons about the top- k most likely proofs, which asymptotically reduces computational
25 cost while providing formal accuracy guarantees relative to the exact instantiation. Scallop thereby
26 generalizes exact probabilistic reasoning and enables easy exploration of a rich space of tradeoffs.
27 This tradeoff mechanism allows to drastically speed up the stochastic training of the involved neural
28 components without sacrificing generalization ability.

29 The main technical contribution of Scallop concerns computing the set of top- k proofs associated with
30 each discrete fact *efficiently*, during the evaluation of a logic program, and *correctly*, by maintaining
31 all and only the top- k proofs. Scallop achieves this goal by formulating the problem in the framework
32 of *provenance* for deductive databases [8]. The framework provides the theory and algorithms for
33 tagging discrete facts derived by a logic program with information—in our case the set of top- k proofs.
34 Concretely, Scallop targets Datalog [1], a syntactic subset of Prolog. Although not Turing-complete,
35 Datalog supports recursion and is expressive enough for a wide variety of applications.

36 Scallop inherits efficient algorithms and optimizations from the databases literature. In contrast,
37 efficiently computing top- k proofs for Prolog is an open problem, to our knowledge. Moreover, the
38 provenance framework enables Scallop to provide correctness guarantees. We leverage the theory of
39 *provenance semirings* [17], which allows us to define how to compute top- k proofs in a compositional
40 manner for each logic operation in Datalog, while ensuring that the computation is correct across

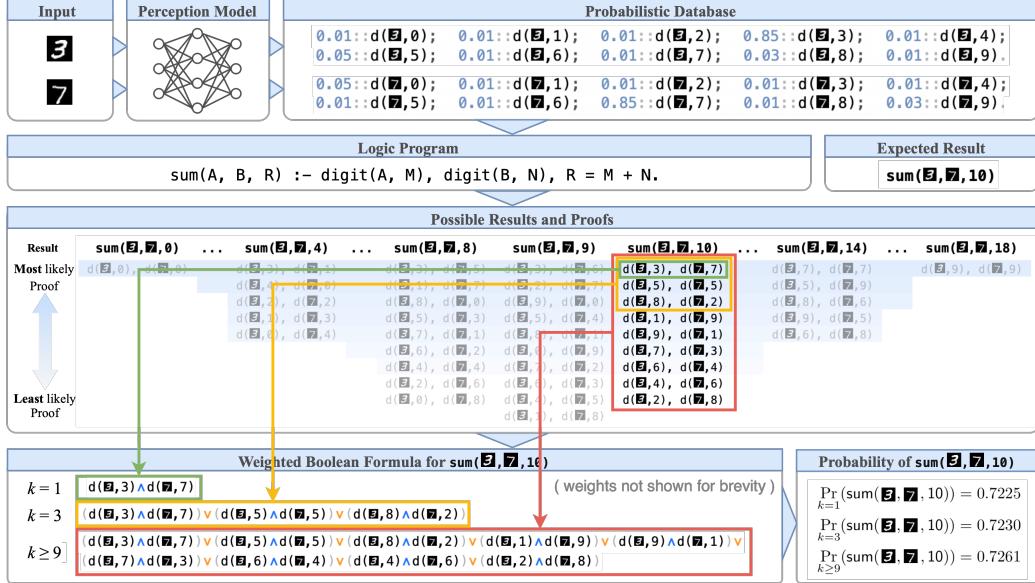


Figure 1: Illustration of our approach on the task $3 + 7 = 10$ using different values of parameter k .

41 arbitrary combinations of these operations. This approach also makes Scallop easy to extend with
42 features such as additional logic operations, probabilistic rules, and foreign functions.

43 We evaluate Scallop on diverse tasks that involve combining perception with reasoning. On a
44 suite of synthetic tasks that involve mathematical and logical reasoning over hand-written digits,
45 Scallop scales significantly better without sacrificing accuracy compared to DeepProbLog [24], a
46 principled neural logic programming approach. We also create and evaluate on a real-world task called
47 VQAR (Visual Question Answering with Reasoning) which augments the VQA task with an external
48 common-sense knowledge base for multi-hop reasoning. The goal is to answer a programmatic
49 question with the correct subset of objects in a real-world image. Scallop takes 92 hours to finish 15
50 training epochs with $k = 10$ and takes only 0.3 seconds on average per training sample. In contrast, a
51 difficult training sample can take DeepProbLog over 100 hours to compute, making it infeasible to
52 train on the whole dataset. Scallop’s differentiable symbolic reasoning pipeline enables it to achieve
53 84.22% test accuracy, outperforming two VQA-tailored neural models based on Neural Module
54 Networks and transformers by 12.42% and 21.66% respectively.

55 In summary, the main contributions of this paper are as follows:

- 56 1. We introduce the notion of top- k proofs which generalizes exact probabilistic reasoning, asymptotically reduces computational cost, and provides relative accuracy guarantees.
57
- 58 2. We develop a framework, Scallop, which introduces a tunable parameter k and efficiently implements the computation of top- k proofs using provenance in Datalog.
59
- 60 3. We empirically evaluate Scallop on synthetic tasks as well as a real-world task, VQA with multi-hop reasoning, and demonstrate that it significantly outperforms baselines.
61

62 2 Illustrative Overview

63 We illustrate our approach using two tasks: a simple task called `sum2` and the real-world VQAR task.

64 **A Simple Task.** The `sum2` task from [24] concerns classifying sums from pairs of hand-written
65 digits, e.g., $3 + 7 = 10$. As depicted in Figure 1, we specify this task using a neural and a symbolic
66 component, following the style of DeepProbLog [24]. The neural component is a perception model
67 that takes in an image of hand-written digit [21] and classifies it into discrete values $\{0, \dots, 9\}$. The
68 symbolic component, on the other hand, is a logic program in Datalog for computing the resulting
69 sum. The interface between the neural and symbolic components is a probabilistic database which
70 associates each candidate output of the perception model with a probability. For instance, the fact
71 $0.85 :: d(3, 3)$ denotes that image 3 is recognized to be the digit 3 with probability 0.85.

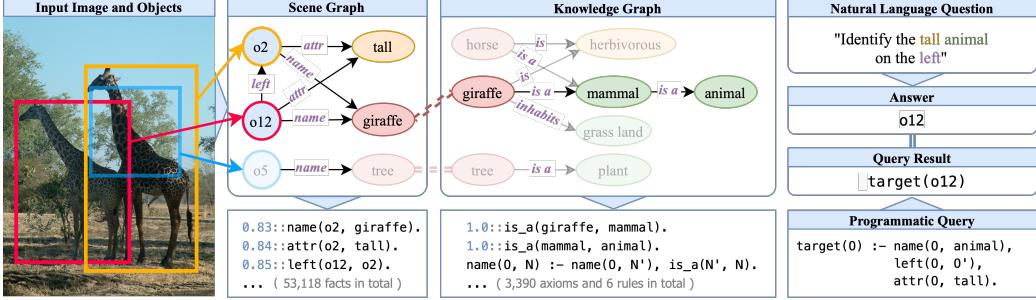


Figure 2: An instance of the VQAR task. The scene graph and knowledge base are shown graphically (above) and in Scallop (below). The question and answer are shown in natural language (above) and in Scallop (below).

72 Evaluating the logic program on the probabilistic database yields a weighted boolean formula for
 73 each possible result of the sum of two digits, i.e., values in the range $\{0, \dots, 18\}$. Each *clause* of
 74 such a formula represents a different *proof* of the corresponding result. For instance, the bottom left
 75 of Figure 1 shows the formula representing all 9 proofs of the ground truth result 10. Each such
 76 formula is input to an off-the-shelf weighted model counting (WMC) solver to yield the probability
 77 of the corresponding result, e.g., $0.7261 :: \text{sum}(\mathbf{3}, \mathbf{7}, 10)$.

78 The scalability of this approach is limited in practice by WMC solving whose complexity is at least
 79 #P-complete [29]. We observe that computing only the top- k most likely proofs bounds the size
 80 of each formula to k clauses, thereby allowing to trade diminishing amounts of accuracy for large
 81 gains in scalability. Moreover, stochastic training of the deep perception models itself can tolerate
 82 noise in data. As we show later in our experiments, the additional noise introduced by the top- k
 83 approximation can be well-compensated for by the stochastic training algorithm.

84 Scallop embodies this insight by introducing a parameter k which can be task-dependent, and even
 85 for a particular task, tuned differently for learning and inference. A higher k leads to slower inference,
 86 but accelerates the convergence of learning, especially for complex or sparse feedback; thus, Scallop
 87 enables to achieve the best of both worlds by employing a higher k during training, and a lower k
 88 thereafter. While Scallop’s inference time is under 0.1 second per task for the `sum2` task regardless
 89 of the choice of k , the difference is much more pronounced for the `sum3` task of adding three digits:
 90 0.05 seconds for $k = 1$ versus 6.15 seconds for $k = 15$.

91 **Visual Question Answering.** We next illustrate applying Scallop to a complex real-world task,
 92 Visual Question Answering (VQA) [4], which is widely studied in the deep learning literature. The
 93 task concerns answering a given question using knowledge from a given image of a scene. Since
 94 we are interested in tasks that combine perception with reasoning, we extend the VQA task with
 95 *multi-hop reasoning* over an external common-sense knowledge base. The resulting task, which we
 96 call *VQAR*, improves upon the VQA task in two important ways: it generalizes the VQA task by
 97 allowing questions that require external knowledge, and it allows to precisely control the reasoning
 98 complexity through the number of hops needed to answer them. [1] We thereby develop a new dataset
 99 consisting of real-world images of scenes and object identification questions that necessitate varying
 100 hops of reasoning in a fixed external knowledge base.

101 It is natural to express the VQAR task using a combination of neural and symbolic modules akin to
 102 the `sum2` task. As Figure 2 illustrates, these modules are more complex, reflecting the real-world
 103 nature of this task. The neural module is a perception model that takes the object feature vectors
 104 (extracted by pre-trained vision models) and outputs a scene graph comprising the predicted name
 105 and attribute distributions of each object, and relationships between the objects—all of which are
 106 uniformly represented as a probabilistic database. For instance, the tuple $0.83 :: \text{name}(o12, \text{giraffe})$
 107 denotes that name of object $o12$ is classified as *giraffe* with probability 0.83.

108 Likewise, the symbolic module uniformly represents both the logic representation of the question
 109 and the external knowledge base as a logic program in Datalog.² Evaluating the program on the
 110 probabilistic database yields the answer, e.g., $\text{target}(o12)$. The example in Figure 2 highlights the

¹In contrast, prior works such as the GQA dataset [19] are limited to varying the reasoning complexity in the question alone, which renders the question unwieldy.

²We presume that the input question is in programmatic form because existing models for semantic parsing achieve high accuracy in translating from natural language text to programmatic form [7].

(Constant)	c	(Probability)	p
(Variable)	V	(Prob. Input Fact)	f $p :: \bar{f} \in \mathcal{F}$
(Term)	$t \in V \mid c$	(Disjunction)	$j \in f_1; \dots; f_n \in \mathcal{J}$
(Predicate)	a	(Query)	$\mathcal{Q} \in \alpha$
(Atom)	$\alpha \in a(t_1, \dots, t_n)$	(Query Result)	$q \in g$
(Fact)	$g \in a(c_1, \dots, c_n) \in \mathcal{G}$	(Program)	$\mathcal{P} \in (\bar{\mathcal{F}}, \mathcal{R}, \mathcal{Q})$
(Input Fact)	$\bar{f} \in g \in \bar{\mathcal{F}}$	(Prob. Program)	$\bar{\mathcal{P}} \in (\mathcal{F}, \mathcal{R}, \mathcal{J}, \mathcal{Q})$
(Rule)	$r \in \alpha := \alpha_1, \dots, \alpha_m \in \mathcal{R}$		

Figure 3: Abstract syntax of probabilistic Datalog programs.

111 need for external knowledge: although the question refers to the concept of an “animal” that is missing
 112 in the scene graph, Scallop is able to derive the conclusion $\text{name}(o_{12}, \text{animal})$ without changing the
 113 perception model. The derivation involves two-hop reasoning—two applications of the recursive rule
 114 $\text{name}(O, N) := \text{name}(O, N'), \text{is_a}(N', N)$ to facts from the scene and knowledge graphs:

$$\frac{\text{name}(o_{12}, \text{giraffe}) \quad \text{is_a}(\text{giraffe}, \text{mammal})}{\text{name}(o_{12}, \text{mammal}) \quad \text{is_a}(\text{mammal}, \text{animal})} \\ \text{name}(o_{12}, \text{animal})$$

115 While more sophisticated models can learn the representation of concepts such as animal from a large
 116 corpus, relying on such pretrained representation sacrifices the benefits of symbolic reasoning, such
 117 as interpretability, data efficiency, and generalization to unseen concepts.

118 3 Background

119 We recap Datalog, the logic programming language that underlies Scallop, and present its probabilistic
 120 extensions that we leverage for inference and training tasks.

121 **Syntax of Datalog.** As shown in Figure 3, a Datalog program $\bar{\mathcal{P}}$ consists of a set of input facts $\bar{\mathcal{F}}$, a
 122 set of rules \mathcal{R} , and a query \mathcal{Q} . The building block is an atom $a(t_1, \dots, t_n)$ which consists of an n -ary
 123 predicate a and a list of terms t_1, \dots, t_n as arguments. A fact g is an atom which all the argument
 124 terms are constants; it may be an input fact (EDB) or a derived fact (IDB). Datalog rules are of the
 125 form $\alpha := \alpha_1, \dots, \alpha_m$, meaning that atom α in the head is true if all atoms α_i in the body are true.
 126 Multiple rules sharing a single head predicate denote disjunction (or union).

127 **Semantics of Datalog.** Datalog programs can be executed using a bottom-up evaluation strategy.
 128 Starting from the input facts $\bar{\mathcal{F}}$, we repeatedly apply the rules \mathcal{R} in any order to derive new facts
 129 until a fixed point is reached. Upon completion, we obtain all the output facts q of the query \mathcal{Q} .
 130 For example, with $\bar{\mathcal{F}} = \{\text{left}(o_1, o_2), \text{below}(o_2, o_3)\}$ and $\mathcal{Q} = \text{left}(o_1, O)$, the execution of program
 131 $(\bar{\mathcal{F}}, \emptyset, \mathcal{Q})$ produces $\{\text{left}(o_1, o_2)\}$. We denote the execution result as $\text{Exec}(\bar{\mathcal{P}}) = \{q_i\}_{i=1}^n$.

132 **Probabilistic Extensions.** To handle uncertain data, we introduce two probabilistic extensions to
 133 Datalog, which are inspired by pD [5] and ProbLog [12]. First, we specify probabilistic input facts
 134 f by associating a probability p with \bar{f} , declaring that $\Pr(f) = p$. Deterministic input facts have
 probability 1.0. Secondly, we allow disjunctions \mathcal{J} among probabilistic input facts, denoted by
 $f_1; \dots; f_m$. For example, the disjunction

$$0.01 :: \text{digit}(\mathbf{3}, 0); \dots; 0.82 :: \text{digit}(\mathbf{3}, 3); \dots; 0.06 :: \text{digit}(\mathbf{3}, 9).$$

135 states that the digit $\mathbf{3}$ is recognized to be 0 to 9 with their respective probabilities, but cannot be
 136 more than one simultaneously. \mathcal{F} and \mathcal{J} form a *probabilistic database*. By combining the \mathcal{F} , \mathcal{J} with
 137 \mathcal{R} and \mathcal{Q} , we obtain a probabilistic Datalog program \mathcal{P} .

138 **Probability Calculation.** Unlike discrete Datalog, which provides definite answers to queries, we
 139 wish to compute the *success probability* of each query result q : $\text{Exec}(\mathcal{P}) = \{(q_i, \Pr(q_i))\}_{i=1}^n$. To
 140 compute success probabilities, we first define a *proof* of any fact g as a minimal set of (probabilistic)
 141 input facts f that can derive g . We denote a proof as $F \in \wp(\mathcal{F})$ where \wp denotes power set. Since
 142 a fact g may be explained by multiple proofs, we use S_g to denote the complete set of proofs of g .
 Given the set of proofs S_q for a query result q , the success probability $\Pr(q)$ is simply the likelihood
 143 of S_q , denoted $\Pr(S_q)$, which can be computed using *Weighted Model Counting* (WMC) [20].

144 4 Framework

145 The Scallop framework aims to solve two sub-problems:

1. **Inference** (Section 4.1): Given a probabilistic Datalog program $\mathcal{P} = (\mathcal{F}, \mathcal{R}, \mathcal{J}, \mathcal{Q})$, efficiently
 146 compute each query result q_i with its set of proofs S_{q_i} .

$$\begin{array}{c}
\frac{\begin{array}{c} f_1 \\ \text{name(o}_{12}\text{, giraffe)} \quad f_2 \\ S_{f_1} = \{\{f_1\}\} \quad S_{f_2} = \{\{f_2\}\} \end{array}}{\frac{\begin{array}{c} g \\ \text{name(o}_{12}\text{, mammal)} \\ S_g = \{\{f_1, f_2\}\} \end{array}}{\text{[AND]}}} \quad \frac{\begin{array}{c} f_3 \\ \text{is_a(mammal, animal)} \\ S_{f_3} = \{\{f_3\}\} \end{array}}{\frac{\begin{array}{c} q : \text{name(o}_{12}\text{, animal)} \\ S_q = \{\{f_1, f_2, f_3\}\} \end{array}}{\text{[AND]}}}} \\
\end{array}$$

Figure 4: Proof constr. with conjunction.

$$\frac{\begin{array}{c} f_1 : \text{name(o}_3\text{, giraffe)} \\ S_{f_1} = \{\{f_1\}\} \end{array}}{\frac{\begin{array}{c} f_2 : \text{name(o}_3\text{, tiger)} \\ S_{f_2} = \{\{f_2\}\} \end{array}}{\frac{\begin{array}{c} q : \text{target(o}_3\text{)} \\ S_q = \{\{f_1\}, \{f_2\}\} \end{array}}{\text{[OR]}}}}} \\$$

Figure 5: Proof constr. with disjunction.

146 2. **Learning** (Section 4.2): Given a neural symbolic reasoning dataset \mathcal{D} and a loss function \mathcal{L} ,
147 learn a perception model \mathcal{M}_θ which, for each $(x, y) \in \mathcal{D}$, transforms x into a probabilistic
148 database captured by Datalog program \mathcal{P}_θ^x . We aim to minimize the following objective: $J(\theta) =$
149 $\frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} \mathcal{L}(\text{Exec}(\mathcal{P}_\theta^x), y)$.

150 4.1 Inference

151 **Proof Construction.** The goal of our proof construction is to construct the set of proofs S_q for every
152 query result q . We can efficiently compute S_q during the bottom-up execution of the Datalog program.
153 We initially tag each input fact $f \in \mathcal{F}$ with $S_f = \{\{f\}\}$ and propagate proofs during execution from
154 known facts to newly derived facts.

155 We illustrate proof propagation during conjunction in Figure 4. When g is derived from a conjunction
156 on f_1 and f_2 , we combine the sets of proofs S_{f_1} and S_{f_2} to produce S_g . The resulting S_g contains a
157 single proof $\{f_1, f_2\}$, as both f_1 and f_2 must be true for g to be true. More formally, we define a
158 binary operation \otimes corresponding to conjunction. Given two sets of proofs S_1 and S_2 , we have

$$S_1 \otimes S_2 = \{F \mid F = F_1 \cup F_2, (F_1, F_2) \in S_1 \times S_2, F \text{ contains no disjunction conflict}\}. \quad (1)$$

159 We next illustrate proof propagation during disjunction in Figure 5. Consider a VQAR instance in
160 which the query concerns identifying a target object that is either a giraffe or a tiger. S_q contains two
161 separate proofs, one containing only f_1 and the other containing only f_2 , as each can individually
162 explain q . We thereby define a binary operation \oplus corresponding to disjunction, as set union:

$$S_1 \oplus S_2 = S_1 \cup S_2. \quad (2)$$

163 Equipped with \oplus and \otimes , we can show that the collection of sets of proofs $\mathcal{S} = \wp(\wp(\mathcal{F}))$ forms a
164 semiring, which we call the *proof semiring*. Following [17], every derivable fact g can be annotated
165 with a corresponding algebraic formula representing the bottom-up construction of S_g . Since
166 the proof semiring is both commutative and distributive, we show in Appendix A.1 that $S_q =$
167 $\bigoplus_F \text{derives } q \left(\bigotimes_{f \in F} S_f \right)$.

168 However, the complexity of S_q renders the computation infeasible. In principle, we have $|S_q| =$
169 $\mathcal{O}(2^{|\mathcal{F}|})$, showing that $|S_q|$ grows exponentially with the amount of input facts. The actual version
170 of our example shown in Figure 2 generates 2,619 proofs in total for all query results, and takes 14
171 minutes to execute. This scalability issue is further exacerbated when the system is used in a learning
172 setting, where we need to execute millions of such programs.

173 **Top- k Proof Construction.** The probabilistic nature of our problem setting opens up room for
174 approximation. A key observation is that, when the inference system is used in a learning setting,
175 the probability of a ground truth fact should significantly outweigh other facts, forming a skewed
176 distribution. We can exploit this property by only including the “most likely” proofs in S_q , with the
177 likelihood of a proof F defined by $\Pr(F) = \prod_{f \in F} \Pr(f)$.

178 We thereby introduce a *top- k proof inference* algorithm. With a user-specified hyper-parameter $k \geq 1$,
179 we perform top- k filtering at each step of the proof construction. We define two new operations, $\otimes^{(k)}$
180 for conjunction, and $\oplus^{(k)}$ for disjunction:

$$S_1 \otimes^{(k)} S_2 = \text{Top}_k(S_1 \otimes S_2), \quad S_1 \oplus^{(k)} S_2 = \text{Top}_k(S_1 \oplus S_2). \quad (3)$$

181 Intuitively, whenever \otimes or \oplus is performed, we rank proofs by their likelihood and preserve only the
182 top- k proofs. This allows us to discard the vast majority of proofs and thus make inference tractable.
183 An example run-through of *top-3 natural join* ($\otimes^{(3)}$) is depicted in Figure 6, where we perform a
184 normal \otimes operation followed by a top-3 filtering.

185 As before, we construct a *top- k proof semiring* (Appendix A.2), with which we can express the
186 resulting approximated *beam of proofs* $\tilde{S}_q = \bigoplus_F \text{derives } q \left(\bigotimes_{f \in F} S_f \right)$. Note that the size of \tilde{S}_q

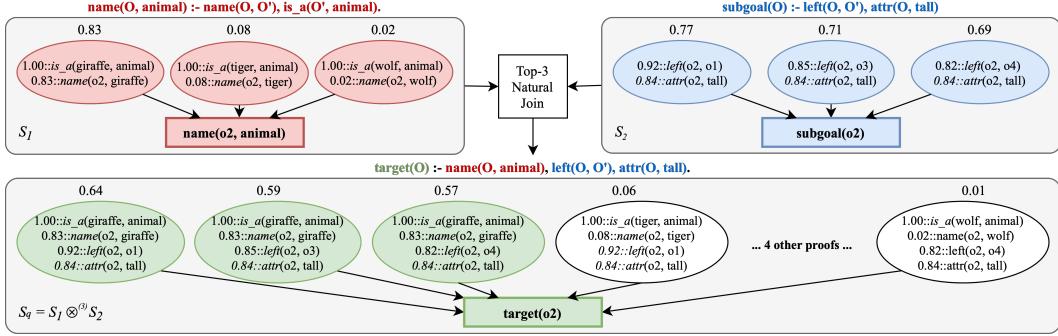


Figure 6: Illustration of top- k natural join using $k = 3$. Each ellipse represents a proof of the fact shown in the box. Given the top 3 proofs for each of “ $\text{name}(o_2, \text{animal})$ ” and “ $\text{subgoal}(o_2)$ ”, we wish to derive the top 3 proofs for their conjunction, “ $\text{target}(o_2)$ ”. The join yields 9 possible proofs. After computing the likelihood for each of the 9 proofs, we keep the top 3 most likely ones (green ellipses) and discard the rest (white ellipses).

is bounded by k , $|\tilde{S}_q| = \mathcal{O}(k)$, reducing the exponential complexity of exact inference to a near constant one. As a comparison point, with top-3 proof inference, the full example shown in Figure 2 only generates 39 proofs, taking only 0.5 seconds to execute. Formally, our approximation of the success probability of a given query result q can be written as $\Pr(q) = \Pr(S_q) \approx \Pr(\tilde{S}_q)$.

Discussion. We present some desirable properties of our top- k inference algorithm. The approximation error bound is given by $|\Pr(S_q) - \Pr(\tilde{S}_q)| \leq \sum_{F \in S_q \setminus \tilde{S}_q} \Pr(F)$, and we can tune k to control the trade-off between scalability and accuracy. Furthermore, if no disjunctions are specified ($\mathcal{J} = \emptyset$), then we have $\tilde{S}_q = \text{Top}_k(S_q)$, that is, the beam of proofs \tilde{S}_q contains the global top- k proofs. The theorems and proofs are provided in Appendix A.3.

We also note that our top- k inference algorithm is reminiscent of beam search. Both methods are iterative and explore only the top- k elements at each step. However, there are two major differences that distinguish us from beam search. First, while beam search is heuristic, our algorithm is backed by Datalog semantics and the provenance semirings framework for its correctness. We also present formal guarantees on its approximation error bound. Secondly, our algorithm operates over the beam of proofs \tilde{S}_q for each derived fact q , while beam search is usually performed to search for an output.

4.2 Learning

At a high level, we want to train a perception model \mathcal{M}_θ that takes in an input x and produces a probabilistic database $(\mathcal{F}, \mathcal{J})$, captured by program \mathcal{P} , such that after execution, can derive the ground truth y as the output. Note that the probability of the input facts in the probabilistic database is generated by the perception model \mathcal{M}_θ . Therefore each input probability $p_i = \Pr(f_i)$ is also associated with their gradients $\nabla_{\Pr(f_i)}$ with respect to the model parameters θ .

To back-propagate the gradients through the inference process, similar to DeepProbLog [24], Scallop adopts a *gradient semiring* augmented WMC procedure, for which we use *Sentential Decision Diagram* (SDD) [10]. The beam of proofs \tilde{S}_q will be transformed into a weighted Conjunctive Normal Form (CNF) formula, where for each variable, f_i , we attach the dual number $(\Pr(f_i), \nabla_{\Pr(f_i)})$ as its weight. As a result, the associated differentiable probability of each query result q_i will be $(\Pr(q_i), \nabla_{\Pr(q_i)})$, as computed by WMC. With everything above, we define the execution of our probabilistic Datalog program as

$$\hat{y} = \text{Exec}(\mathcal{P}) = \{(q_i, (\Pr(q_i), \nabla_{\Pr(q_i)}))\}_{i=1}^n. \quad (4)$$

The results of the execution \hat{y} , along with the ground truth y is passed to the given loss function \mathcal{L} . Lastly, the loss is back-propagated to update θ , the parameters of the perception model \mathcal{M}_θ .

For example, the ground truth label y for the task $\text{sum}(\mathbf{3}, \mathbf{7}, \text{R})$ is a binary vector of dimension 19, conceptually representing the set:

$$\{0.0 :: \text{sum}(\mathbf{3}, \mathbf{7}, 0), \dots, 1.0 :: \text{sum}(\mathbf{3}, \mathbf{7}, 10), \dots, 0.0 :: \text{sum}(\mathbf{3}, \mathbf{7}, 18)\}.$$

and the predicted \hat{y} is a set of the 19 results associated with their predicted probabilities, represented as a probability vector of dimension 19. In our experimental setup, we apply the binary cross entropy loss function on the two vectors. In practice, however, the loss function is fully customizable.

Task	Goal Predicate	#Out	Max #Proofs	Scallop				DPL
				$k=1$	$k=3$	$k=5$	$k=10$	
T1	sum2(3 , 7 , 10)	19	10	97.46%	96.90%	96.67%	96.29%	96.82%
T2	sum3(3 , 7 , 5 , 15)	28	75	95.31%	95.43%	95.76%	95.76%	95.56%
T3	sum4(3 , 7 , 5 , 2 , 17)	37	670	47.11%	95.47%	95.31%	95.07%	—
T4	sort2(3 , 7 , 0, 1)	2	55	80.43%	91.55%	91.75%	95.49%	98.04%
T5	sort3(7 , 2 , 3 , 1, 2, 0)	6	220	70.34%	93.20%	96.15%	97.09%	95.50%
T6	sort4(7 , 3 , 5 , 2 , 3, 1, 2, 0)	24	715	68.67%	87.90%	92.02%	91.87%	89.96%

Table 1: Testing accuracy of Scallop and DeepProbLog (DPL) on a suite of 6 synthetic tasks. All numbers except $k = 1$ have a standard deviation of < 2%.

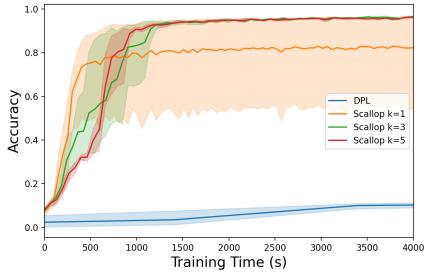


Figure 7: Training runtime (in seconds) vs. validation accuracy for task **T2** (sum3).

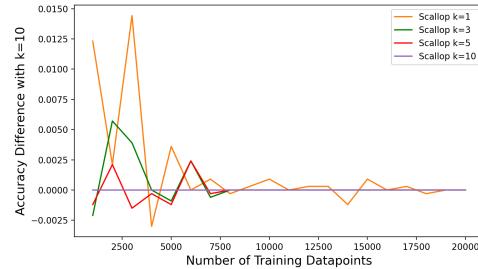


Figure 8: Difference in accuracy of varying k_{test} compared to $k_{\text{test}} = 10$ for task **T2** (sum3).

5 Evaluation

We evaluate Scallop on a suite of synthetic tasks and VQAR. All experiments are conducted on a machine with two 20-core Intel Xeon CPUs, four GeForce RTX 2080 Ti GPUs, and 768 GB RAM. Experimental details such as hyperparameter selection and dataset splits are provided in Appendix C and implementation details of the Scallop framework are explained in Appendix D.

5.1 Synthetic Tasks

We extend the synthetic tasks from DeepProbLog (DPL) to demonstrate that (1) Scallop is much more scalable, (2) Scallop does not sacrifice accuracy, and (3) how different levels of reasoning granularity during training and testing phases can affect model performance.

Table I shows 6 synthetic tasks and their corresponding sample goal predicates. Each task takes as input multiple MNIST [21] images and requires performing simple arithmetic (T1-T3) or sorting (T4-T6) over digits depicted in the given images. The difficulty of each task is reflected by third and fourth columns, which show the size of the output space and the maximum number of proofs per output, respectively. Our goal is to train a digit classifier end-to-end with the combined perception + reasoning pipeline. We elaborate on individual tasks further in Appendix E.

Accuracy. We show accuracy comparison with DPL in Table I. All models are trained under the same learning setting. Scallop is able to achieve on par accuracy as DPL, despite using far fewer proofs. It also shows that in general, larger k implies better accuracy. Note that we are unable to collect result for DPL on T3, as DPL takes 24 hours only to complete 100 out of the 15,000 training samples. In contrast, Scallop with $k = 3$ finishes 5 epochs (75,000 training samples) within 4 hours.

Runtime vs. Accuracy. We next evaluate the tradeoff between the training runtime vs. testing accuracy in Scallop. Figure 7 shows the results for the sum3 task. With $k = 1$, Scallop learns the fastest in the beginning, but it has high variance and potential of failing to converge to an optimal solution. On the other hand, with $k = 5$, it has much less variance and converges the fastest despite being slower in the beginning. We compare with DPL trained under the same setting. It achieves the same accuracy (95.56%) at the end of the 3rd epoch, but due to its long runtime (14 hours), we omit showing the whole curve in this figure.

Decoupling Reasoning Granularity. Scallop enables using different k during training and testing phases. The key idea is that a larger k will help faster convergence in training, whereas a smaller k suffices during testing since less probable proofs have minimal impact on the reasoning result. In Figure 8, we fix a $k_{\text{train}} = 10$ on the sum3 task. Taking accuracy with $k_{\text{test}} = 10$ as a baseline, we compute the difference in testing accuracy on $k_{\text{test}} \in \{1, 3, 5\}$. The figure shows that as the training

Test Dataset	LXMERT	NMN	Scallop
1000 C2	66.75%	79.32%	85.17%
1000 C3	61.69%	61.98%	82.82%
1000 C4	63.82%	71.17%	83.25%
1000 C5	64.05%	74.62%	85.53%
1000 C6	56.51%	72.04%	84.30%
5000 C_{all}	62.56%	71.80%	84.22%

Table 2: Testing accuracy (in Recall@5) of Scallop, NMNs, and LXMERT on VQAR dataset.

252 progresses, the difference converges to 0%. This suggests we can tune k_{train} and k_{test} individually for
253 better training as well as faster test time inference.

254 5.2 Visual Question Answering

255 We next evaluate Scallop on the VQAR task described in Section 2. Besides DPL, we compare with
256 two neural methods: Neural Module Network (NMN) and LXMERT, a transformer based approach.

257 **Dataset.** The VQAR dataset contains (a) 80,178 images, (b) object feature vectors + bounding boxes,
258 (c) scene graphs with 500 object names, 609 attributes, and 229 relationships, (d) a shared knowledge
259 graph with 6 rules and 3K knowledge triplets, and (e) 4M programmatic queries and answer pairs.
260 The images and scene graphs are from the GQA [19] dataset and the knowledge graph is from the
261 CRIC [16] dataset. The object feature vectors and bounding boxes are then obtained by passing the
262 images through pre-trained fixed-weight Mask RCNN and ResNet models. Using random walk on
263 combined scene graph and external knowledge graph, we generate object identification questions
264 in the form of programmatic queries. We further categorize these queries into different levels of
265 difficulty by the number of occurring clauses from C2 to C6, where C2 is the simplest and C6 is the
266 hardest. Further details of this dataset are provided in Appendix B.

267 We formulate VQAR as a multi-label classification task. For each datapoint (x, y) in our VQAR
268 dataset, the input x consists of (a) the shared knowledge graph KG , (b) a programmatic query, and
269 (c) the object feature vectors and bounding boxes. The ground truth y is the set of objects that the
270 given programmatic query identifies. All of our evaluated models share this same set of input and
271 output (except LXMERT, which takes in natural language questions instead of programmatic queries).
272 The accuracy is measured by Recall@5.

273 **Setup of Scallop.** We use a perception module consisting of three MLP-classifiers, $\mathcal{M}_\theta =$
274 $(\mathcal{M}_\theta^n, \mathcal{M}_\theta^a, \mathcal{M}_\theta^r)$, which predict names, attributes, and relations respectively. All predictions are
275 transformed into probabilistic facts in a database. The outputs of \mathcal{M}_θ^n form disjunctions because
276 each object has only one name. With KG as part of the probabilistic database, we perform Datalog
277 execution on the given programmatic query to obtain the set of identified objects. Note that the
278 *entire* knowledge graph is used in every Datalog execution. We use binary cross entropy as our loss
279 function to compare the predicted set of objects and the ground truth set. The goal is to train the three
280 classifiers in Scallop end-to-end, and identify the correct objects according to the question.

281 **Baseline 1: DeepProbLog.** It is prohibitively slow to train with DPL from scratch—a regular
282 training sample from C6 can take DPL more than 100 hours to run. Therefore, instead of training
283 with DPL, we use the perception model \mathcal{M}_θ trained with Scallop to test DPL’s inference capability.
284 With 10 seconds timeout, DPL times out on 68.66% of the testing samples, while Scallop finishes all
285 with an average running time under 0.3 seconds per sample.

286 **Baseline 2: Neural Module Network.** We compare against RVC [16], a Neural Module Network
287 approach for VQA with external common-sense knowledge. This method first pretrains a TransE
288 embedding [6] for the knowledge graph. Then, to mimic the reasoning process, it trains a set of
289 neural modules that perform knowledge retrieval, scene graph traversal, and logical operations. The
290 modules are assembled according to the programmatic query and can leverage object-based features.

291 **Baseline 3: LXMERT.** We also compare to LXMERT [32], a recent transformer based approach that
292 emphasizes its transfer learning ability. LXMERT takes in a natural language question corresponding
293 to the given programmatic query. Similar to other baselines, the object features and bounding boxes
294 are taken as input. Since this model cannot explicitly use a knowledge base, we leverage the implicit
295 relations learned through pre-training over a variety of image-language tasks: MS COCO [23], Visual
296 Genome [4], and GQA [19]. Finally, we fine-tune LXMERT on our VQAR training samples.

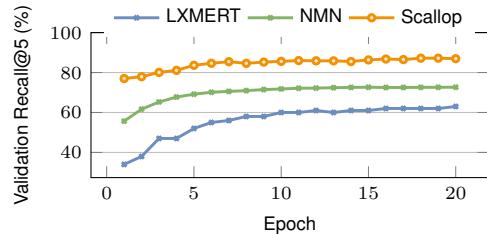


Figure 9: Results of training on 50K C_{all} tasks and testing on 5000 tasks of different clause lengths.

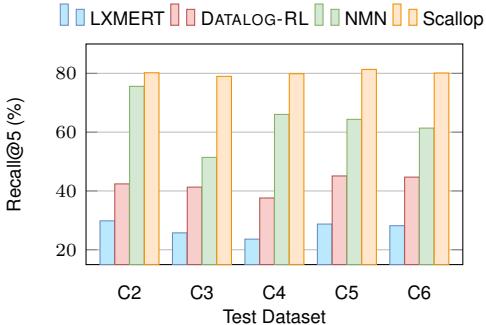


Figure 10: Generalizability to harder questions when trained on 10K C2.

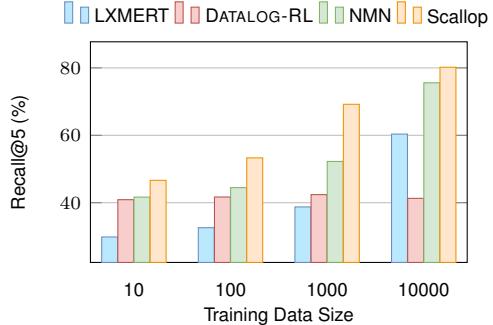


Figure 11: Data efficiency given training data size from 10 to 10,000 C2.

297 **Ablation Study: Datalog Reinforcement Learning (DATALOG-RL).** In this study, we remove the
 298 differentiability in Scallop’s learning pipeline. Instead, we sample a discrete scene graph, run it
 299 through the standard Datalog execution, and use the overlap in predicted objects as a reward to
 300 estimate the gradient using REINFORCE [36]. This method does not scale with the training dataset
 301 of 50K tasks, so we only perform the generalizability experiments (Figure 10).

302 **Results.** Table 2 and Figure 9 compares the performance of Scallop, NMN, and LXMERT based on
 303 50K training tasks. Scallop significantly outperforms both in terms of accuracy and data efficiency.
 304 Figure 10 shows that Scallop generalizes to answer more difficult questions (1K from each of C2-C6)
 305 even when trained on only the easiest ones (10K C2). Figure 11, on the other hand, shows the testing
 306 accuracy (on 1K C2) when trained on varying dataset sizes (10, 100, 1000, and 10,000 C2). We
 307 observe that Scallop has the best data efficiency. Finally, with DATALOG-RL we observe that the
 308 addition of differentiable reasoning is crucial to Scallop’s learning performance.

309 6 Related Work

310 Neural symbolic methodology aims to disentangle low-level perception from high-level reasoning
 311 systematically. Generically speaking, there are three classes of the neural symbolic method. (1) Logic
 312 regularization term. Whenever the network fails to obey the logic constraint, it will receive a penalty
 313 [30, 37]. (2) Soft logic program execution. The primitive operations in a logic program are mapped
 314 to differentiable mathematical operations or neural components [14, 28]. (3) Proof-guided probability
 315 calculation. Approaches like exact probability calculation and abductive reasoning first execute the
 316 logic program and then map the generated proof constructs into differentiable expressions [9, 22, 24].

317 Using logic constraints as regularization terms can scale, but does not guarantee the reasoning
 318 correctness. Substituting logic reasoning steps by differentiable components fails to preserve the
 319 original semantics of logic reasoning. Exact probability calculation, on the other hand, maintains
 320 the purity of the logic reasoning pipeline, but has significant scalability limitation. Most application-
 321 specific neural symbolic approaches fall in categories (1) and (2) due to their high-efficiency demand.

322 In the VQA task, a natural language question can be parsed into a logic program that takes in the
 323 image and generates the answer. There are a few studies on how to execute such logic programs
 324 softly [2, 3, 25, 38]. In a neural-guided program synthesis task, given a set of input-out examples, a
 325 neural module is trained to guide the program search process. As programs with complex syntactic
 326 components are non-trivial to transform into differentiable operations, many approaches utilize
 327 reinforcement learning to provide feedback [5, 18, 31, 33, 34].

328 7 Conclusion and Future Work

329 We proposed Scallop, a framework for scaling differentiable reasoning based on Datalog, motivated by
 330 real-world applications that necessitate combining perception and reasoning. The key idea underlying
 331 Scallop is to relax exact probabilistic reasoning via a tunable parameter that specifies the level of
 332 reasoning granularity. We demonstrated the effectiveness of Scallop on diverse tasks including a
 333 newly created Visual Question Answering benchmark that requires multi-hop reasoning. In future, we
 334 plan to develop expressive extensions to Scallop, target more challenging neuro-symbolic applications,
 335 and optimize the end-to-end pipeline on modern hardware.

336 **References**

- 337 [1] ABITEBOUL, S., HULL, R., AND VIANU, V. *Foundations of Databases: The Logical Level*,
338 1st ed. Pearson, 1994.
- 339 [2] AMIZADEH, S., PALANGI, H., POLOZOV, O., HUANG, Y., AND KOISHIDA, K. Neuro-
340 symbolic visual reasoning: Disentangling" visual" from" reasoning". *arXiv preprint arXiv:2006.11524* 2 (2020).
- 341 [3] ANDREAS, J., ROHRBACH, M., DARRELL, T., AND KLEIN, D. Neural module networks.
342 In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016),
343 pp. 39–48.
- 344 [4] ANTOL, S., AGRAWAL, A., LU, J., MITCHELL, M., BATRA, D., ZITNICK, C. L., AND
345 PARIKH, D. Vqa: Visual question answering. In *Proceedings of the IEEE international
346 conference on computer vision* (2015), pp. 2425–2433.
- 347 [5] BALOG, M., GAUNT, A. L., BROCKSCHMIDT, M., NOWOZIN, S., AND TARLOW, D. Deep-
348 coder: Learning to write programs. *CoRR abs/1611.01989* (2016).
- 349 [6] BORDES, A., USUNIER, N., GARCIA-DURAN, A., WESTON, J., AND YAKHNENKO, O.
350 Translating embeddings for modeling multi-relational data. In *Neural Information Processing
351 Systems (NIPS)* (2013), pp. 1–9.
- 352 [7] CHEN, W., GAN, Z., LI, L., CHENG, Y., WANG, W., AND LIU, J. Meta module network
353 for compositional visual reasoning. In *Proceedings of the IEEE/CVF Winter Conference on
354 Applications of Computer Vision* (2021), pp. 655–664.
- 355 [8] CHENEY, J., CHITICARIU, L., AND TAN, W.-C. Provenance in databases: Why, how, and
356 where. *Foundations and Trends in Databases* 1, 4 (Apr. 2009).
- 357 [9] DAI, W.-Z., XU, Q., YU, Y., AND ZHOU, Z.-H. Bridging machine learning and logical
358 reasoning by abductive learning. In *NeurIPS 2019* (2019).
- 359 [10] DARWICHE, A. Sdd: A new canonical representation of propositional knowledge bases. In
360 *Proceedings of the Twenty-Second International Joint Conference on Artificial Intelligence -
361 Volume Volume Two* (2011), IJCAI'11, AAAI Press, p. 819–826.
- 362 [11] D'AVILA GARCEZ, A., GORI, M., LAMB, L. C., SERAFINI, L., SPRANGER, M., AND TRAN,
363 S. N. Neural-symbolic computing: An effective methodology for principled integration of
364 machine learning and reasoning, 2019.
- 365 [12] DE RAEDT, L., KIMMIG, A., AND TOIVONEN, H. Problog: A probabilistic prolog and its
366 application in link discovery. pp. 2462–2467.
- 367 [13] DEUTCH, D., GILAD, A., AND MOSKOVITCH, Y. Efficient provenance tracking for datalog
368 using top-k queries. *The VLDB Journal* 27 (2018), 245–269.
- 369 [14] EVANS, R., AND GREFENSTETTE, E. Learning explanatory rules from noisy data. *Journal of
370 Artificial Intelligence Research* 61 (2018), 1–64.
- 371 [15] FUHR, N. Probabilistic datalog—a logic for powerful retrieval methods. In *Proceedings
372 of the 18th Annual International ACM SIGIR Conference on Research and Development in
373 Information Retrieval* (New York, NY, USA, 1995), SIGIR '95, Association for Computing
374 Machinery, p. 282–290.
- 375 [16] GAO, D., WANG, R., SHAN, S., AND CHEN, X. From two graphs to n questions: A vqa dataset
376 for compositional reasoning on vision and commonsense. *arXiv preprint arXiv:1908.02962*
377 (2019).
- 378 [17] GREEN, T. J., KARVOUNARAKIS, G., AND TANNEN, V. Provenance semirings. In *Proceedings
379 of the ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems (PODS)*
380 (2007).
- 381

- 382 [18] HUANG, J., SMITH, C., BASTANI, O., SINGH, R., ALBARGHOUTHI, A., AND NAIK, M.
 383 Generating programmatic referring expressions via program synthesis. In *Proceedings of the*
 384 *37th International Conference on Machine Learning* (13–18 Jul 2020), H. D. III and A. Singh,
 385 Eds., vol. 119 of *Proceedings of Machine Learning Research*, PMLR, pp. 4495–4506.
- 386 [19] HUDSON, D. A., AND MANNING, C. D. GQA: a new dataset for compositional question
 387 answering over real-world images. *CoRR abs/1902.09506* (2019).
- 388 [20] KIMMIG, A., DEN BROECK, G. V., AND RAEDT, L. D. Algebraic model counting. *CoRR*
 389 *abs/1211.4475* (2012).
- 390 [21] LECUN, Y., BOTTOU, L., BENGIO, Y., AND HAFFNER, P. Gradient-based learning applied to
 391 document recognition. *Proceedings of the IEEE* 86, 11 (1998), 2278–2324.
- 392 [22] LI, Q., HUANG, S., HONG, Y., CHEN, Y., WU, Y. N., AND ZHU, S.-C. Closed loop neural-
 393 symbolic learning via integrating neural perception, grammar parsing, and symbolic reasoning.
 394 In *International Conference on Machine Learning* (2020), PMLR, pp. 5884–5894.
- 395 [23] LIN, T., MAIRE, M., BELONGIE, S. J., BOURDEV, L. D., GIRSHICK, R. B., HAYS, J.,
 396 PERONA, P., RAMANAN, D., DOLLÁR, P., AND ZITNICK, C. L. Microsoft COCO: common
 397 objects in context. *CoRR abs/1405.0312* (2014).
- 398 [24] MANHAEVE, R., DUMANČIĆ, S., KIMMIG, A., DEMEESTER, T., AND RAEDT, L. D. Deep-
 399 problog: Neural probabilistic logic programming. In *NeurIPS 2018* (2018).
- 400 [25] MAO, J., GAN, C., KOHLI, P., TENENBAUM, J. B., AND WU, J. The neuro-symbolic concept
 401 learner: Interpreting scenes, words, and sentences from natural supervision. *arXiv preprint*
 402 *arXiv:1904.12584* (2019).
- 403 [26] MARINO, K., RASTEGARI, M., FARHADI, A., AND MOTTAQHI, R. Ok-vqa: A visual question
 404 answering benchmark requiring external knowledge. In *Conference on Computer Vision and*
 405 *Pattern Recognition (CVPR)* (2019).
- 406 [27] RAEDT, L. D., MANHAEVE, R., DUMANCIC, S., DEMEESTER, T., AND KIMMIG, A. Neuro-
 407 symbolic = neural + logical + probabilistic. In *International Workshop on Neural-Symbolic*
 408 *Learning and Reasoning* (2019).
- 409 [28] ROCKTÄSCHEL, T., AND RIEDEL, S. End-to-end differentiable proving. *CoRR abs/1705.11040*
 410 (2017).
- 411 [29] SANG, T., BEAME, P., AND KAUTZ, H. A. Performing bayesian inference by weighted model
 412 counting. In *AAAI* (2005), vol. 5, pp. 475–481.
- 413 [30] SERAFINI, L., AND D’AVILA GARCEZ, A. S. Logic tensor networks: Deep learning and
 414 logical reasoning from data and knowledge. *CoRR abs/1606.04422* (2016).
- 415 [31] SI, X., RAGHOTHAMAN, M., HEO, K., AND NAIK, M. Synthesizing datalog programs using
 416 numerical relaxation. *CoRR abs/1906.00163* (2019).
- 417 [32] TAN, H., AND BANSAL, M. Lxmert: Learning cross-modality encoder representations from
 418 transformers. *arXiv preprint arXiv:1908.07490* (2019).
- 419 [33] VIJAYAKUMAR, A. J., MOHTA, A., POLOZOV, O., BATRA, D., JAIN, P., AND GULWANI,
 420 S. Neural-guided deductive search for real-time program synthesis from examples. *CoRR*
 421 *abs/1804.01186* (2018).
- 422 [34] VIJAYAKUMAR, A. J., MOHTA, A., POLOZOV, O., BATRA, D., JAIN, P., AND GULWANI,
 423 S. Neural-guided deductive search for real-time program synthesis from examples. *CoRR*
 424 *abs/1804.01186* (2018).
- 425 [35] WANG, P., WU, Q., SHEN, C., HENGEL, A., AND DICK, A. Explicit knowledge-based
 426 reasoning for visual question answering.
- 427 [36] WILLIAMS, R. J. Simple statistical gradient-following algorithms for connectionist reinforce-
 428 ment learning. *Machine learning* 8, 3-4 (1992), 229–256.

- 429 [37] XU, J., ZHANG, Z., FRIEDMAN, T., LIANG, Y., AND VAN DEN BROECK, G. A semantic loss
430 function for deep learning with symbolic knowledge. In *Proceedings of the 35th International*
431 *Conference on Machine Learning* (10–15 Jul 2018), J. Dy and A. Krause, Eds., vol. 80 of
432 *Proceedings of Machine Learning Research*, PMLR, pp. 5502–5511.
- 433 [38] YI, K., WU, J., GAN, C., TORRALBA, A., KOHLI, P., AND TENENBAUM, J. B. Neural-
434 symbolic vqa: Disentangling reasoning from vision and language understanding. *arXiv preprint*
435 *arXiv:1810.02338* (2018).

436 **Checklist**

- 437 1. For all authors...
- 438 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
439 contributions and scope? [Yes]
- 440 (b) Did you describe the limitations of your work? [Yes] We have future work to improve
441 our method.
- 442 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 443 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
444 them? [Yes]
- 445 2. If you are including theoretical results...
- 446 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 447 (b) Did you include complete proofs of all theoretical results? [Yes] In appendix A
- 448 3. If you ran experiments...
- 449 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
450 mental results (either in the supplemental material or as a URL)? [Yes] In supplemental
451 material
- 452 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
453 were chosen)? [Yes] In supplemental material
- 454 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
455 ments multiple times)? [Yes]
- 456 (d) Did you include the total amount of compute and the type of resources used (e.g., type
457 of GPUs, internal cluster, or cloud provider)? [Yes] At the beginning of evaluation
458 section.
- 459 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 460 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 461 (b) Did you mention the license of the assets? [N/A]
- 462 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
463 In supplemental material.
- 464 (d) Did you discuss whether and how consent was obtained from people whose data you're
465 using/curating? [N/A] We generate synthetic dataset.
- 466 (e) Did you discuss whether the data you are using/curating contains personally identifiable
467 information or offensive content? [N/A] We generate synthetic dataset.
- 468 5. If you used crowdsourcing or conducted research with human subjects...
- 469 (a) Did you include the full text of instructions given to participants and screenshots, if
470 applicable? [N/A]
- 471 (b) Did you describe any potential participant risks, with links to Institutional Review
472 Board (IRB) approvals, if applicable? [N/A]
- 473 (c) Did you include the estimated hourly wage paid to participants and the total amount
474 spent on participant compensation? [N/A]