DOLPHIN: A PROGRAMMABLE FRAMEWORK FOR SCALABLE NEUROSYMBOLIC LEARNING

Anonymous authors

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

023

024

025

026 027

048

052

Paper under double-blind review

ABSTRACT

Neurosymbolic learning has emerged as a promising paradigm to incorporate symbolic reasoning into deep learning models. However, existing frameworks are limited in scalability with respect to both the training data and the complexity of symbolic programs. We propose DOLPHIN, a framework to scale neurosymbolic learning at a fundamental level by mapping both forward chaining and backward gradient propagation in symbolic programs to vectorized computations. For this purpose, DOLPHIN introduces a set of abstractions and primitives directly on top of a high-performance deep learning framework like PyTorch. It thereby enables neurosymbolic programs to be written in a language like Python that is familiar to developers and compile them to computation graphs that are amenable to end-toend differentiation on GPUs. We evaluate DOLPHIN on a suite of 13 benchmarks across 5 tasks that combine deep learning models for text, image, or video processing with symbolic programs that involve multi-hop reasoning, recursion, and black-box functions like Python eval (). DOLPHIN achieves comparable or better accuracy on all benchmarks while taking 0.3%-61.7% of the time (and 23.2%) on average) to train these models on the largest input per task compared to baselines Scallop, ISED, and IndeCateR+, which time out on most of these inputs.

028 1 INTRODUCTION

Deep learning has made great strides in tasks such as image classification, speech recognition, and natural language processing. With the emergence of foundation models like GPT-4 and CLIP, deep learning is increasingly applied to more complex tasks. While such models work well for prediction and generation tasks, they are limited in their ability to perform reasoning required for tasks involving structure, logic, and planning, where symbolic approaches traditionally excel (Kambhampati et al.). Neurosymbolic programming (Chaudhuri et al.) [2021] has emerged as a promising paradigm to incorporate symbolic reasoning into deep learning models, providing the best of both worlds.

Various frameworks have been developed to improve the programmability and accessibility of neurosymbolic applications (Manhaeve et al., 2018; Li et al., 2023; Solko-Breslin et al., 2024). These frameworks support complex symbolic reasoning features like recursion and black-box functions, implement efficient differentiable reasoning algorithms, and provide bindings for deep learning frameworks like PyTorch. However, these frameworks incur significant overhead during training.

Consider a typical workflow of such a framework in Figure (1)(a). We have a supervised learning task with labeled data (x, y), a neural network M_{θ} that processes input x, and a symbolic program $P_{symbolic}$ that takes the network's output r and produces final output y. Existing frameworks, such as Scallop (Li et al., 2023), execute the neural model on GPU but use a separate CPU-based backend (implemented in Rust in Scallop's case) for the symbolic program. Moreover, they introduce interprocess latency in transferring state between the neural and symbolic sub-systems.





(a) Typical neurosymbolic framework (e.g. Scallop).



Figure 1: Comparison of system architectures of neurosymbolic frameworks.

054 Together, these issues hinder the scalability of neurosymbolic learning with respect to problem com-055 *plexity* and *data complexity*. First, the symbolic computation engine must derive a set of all possible 056 results and their associated probabilities in a manner that is differentiable with respect to the net-057 work's parameters θ . As the complexity of the symbolic program increases, the number of possible 058 results and their associated weights also grows exponentially, leading to a combinatorial explosion in the number of required computations. However, the symbolic computations are discrete and not easily parallelizable on modern hardware like GPUs. Second, larger datasets also compound the 060 computational cost of neurosymbolic learning. Deep learning typically addresses this challenge 061 by batching computations across multiple data samples. However, in neurosymbolic learning, the 062 computations may differ across data samples, making it difficult to batch them effectively. 063

- To address these challenges, we need to fundamentally rethink the design of a neurosymbolic framework. One approach is to develop specialized and low-level primitives that scale specific benchmarks but make it time-intensive for developers to write neurosymbolic programs tailored to particular tasks. Alternatively, providing high-level primitives—such as a logic programming language like Scallop (Li et al., 2023) or DeepProbLog (Manhaeve et al., 2018)—simplifies the development of symbolic programs but limits the fine-grained control needed to scale specific applications. Finally, to truly democratize neurosymbolic programming, it is crucial to develop a framework that seamlessly integrates into the everyday deep learning workflows that developers already use.
- In this work, we propose DOLPHIN, a novel framework for scalable neurosymbolic learning. In 072 DOLPHIN, we build three key components that effectively tackle the scalability and programmabil-073 ity challenges described above. First, we develop a general symbolic representation that efficiently 074 captures the relationships between neural network outputs and associated discrete symbols. Sec-075 ond, we introduce a set of primitives to map forward chaining in symbolic programs to vectorized 076 computations over these representations. Third, we develop a set of vectorized provenance semir-077 ings (Green et al., 2007) that are easily pluggable into DOLPHIN and enable to efficiently compute symbolic gradients. As illustrated in Figure **Ib**, these components together allow DOLPHIN to build 079 a computation graph that spans both symbolic and neural operations, is highly parallelizable, and end-to-end differentiable on GPUs. Finally, DOLPHIN is implemented as a library that is integrated 081 with PyTorch, allowing users to easily incorporate it into their existing deep learning pipelines.

We evaluate DOLPHIN on a diverse set of neurosymbolic tasks that involve text, image, video, and multi-modal data, and use rich reasoning features such as recursion and black-box Python functions. Neurosymbolic programs written using DOLPHIN only require 0.3%-61.7% (23.2% on average) of the time to train compared to state-of-the-art baselines including differentiable reasoning frameworks like Scallop, and sampling-based frameworks like ISED and IndeCateR+ while maintaining similar levels of accuracy. We also observe that DOLPHIN efficiently scales to more complex benchmarks and larger datasets whereas the baselines either time out after 10 hours or fail to converge.

- We make the following contributions in this work:
 - We propose DOLPHIN, a novel neurosymbolic programming framework for end-to-end differentiable symbolic reasoning in a scalable manner.
 - We develop novel abstractions to represent symbolic and neural computations and introduce vectorized primitives for neurosymbolic programs.
 - We develop vectorized provenances that can be plugged into DOLPHIN for efficient computation of symbolic gradients on parallelizable hardware such as GPUs.
 - We evaluate DOLPHIN on a diverse range of challenging neurosymbolic tasks across different domains and show that it effectively scales with increasing problem complexity and dataset size.

2 OVERVIEW

091

092

093

094

095

096

098 099

100

We illustrate our approach using the MNIST Sum-N task from (De Smet et al.) 2024). The goal is to train a model that takes as input N images of MNIST digits and returns the sum of the digits represented by the images. During learning, supervision is provided only on the sum instead of the labels of the digits. The difficulty of the problem scales exponentially as there are 10^N states in the input space. Further, there are only 9N + 1 possible labels, resulting in very sparse supervision.

Figure 2a shows the code for this task using DOLPHIN with PyTorch. The neural module is a convolutional neural network (CNN) called MNISTNet. It takes in a batch of image tuples imgs where each sample contains N MNIST images. MNISTNet classifies each image into one of 10



 $\begin{array}{c} & \text{apply} \\ & \text{apply} \\ & \text{+} \end{array} \end{array} \xrightarrow{\text{apply}} \\ & \text{+} \end{array} \qquad \begin{array}{c} \text{apply} \\ \text{concat} \end{array} \xrightarrow{\text{D}_{Path}} \\ & \text{D}_{Path} \end{array} \xrightarrow{\text{D}_{Path}} \\ & \text{False} \\ & \text{false} \end{array} \xrightarrow{\text{apply}} \\ & \text{false} \\ & \text{false} \end{array}$

D'Path

union

fixpoint?

(b) PathFinder.

D_{FinalPath}

D_{Res}

(a) MNIST Sum-N.

135

128 129

130

Figure 3: Computation graphs for two neurosymbolic programs written using DOLPHIN.

classes representing the digits 0-9. The logits produced by MNISTNet, representing probability
 distributions over the digits, are then passed as inputs to the symbolic program. Lines 13-17 depict
 a symbolic program written in Python using DOLPHIN primitives.

139 In order to support training, the symbolic program must reason over all the outputs of the CNN, and 140 return probability distributions over all the possible results (0 to 9N). This involves tracking the 141 probabilities of individual symbols (here, digits or numbers), combinatorially evaluating the results 142 of complex symbolic functions, and calculating the probabilities of each intermediate result, all while tracking their gradients to allow for accurate backpropagation while optimizing the training 143 objective. The batched nature of data in machine learning further complicates these calculations 144 since the probabilities of symbols can be different across samples within the same batch. As a 145 result, writing neurosymbolic programs in native PyTorch is tedious even for simple tasks. 146

To address these issues, DOLPHIN provides primitives that allow programmers to express symbolic programs without worrying about the underlying computations. Lines 13 and 15 of Figure 2a show how the CNN's output can be captured within Distribution objects. Each Distribution associates the digits with the corresponding batched logits produced by the CNN, along with any gradients and associated metadata. Figure 2b shows the internal structure of these objects.

152 The programmer can now express the symbolic program in terms of operations that manipulate 153 Distributions. For instance, in line 16, the apply function is used to perform an operation on two distributions. Here, the apply function takes two Distributions as arguments, along with 154 a lambda function that specifies the addition operation. Under the hood, apply combinatorially 155 explores all the possible sums of the symbols from D_res and D_i and calculates their associated 156 probabilities using an appropriate provenance. The result of apply is a new Distribution over 157 the calculated sums, and is stored back into D_res. This is repeated iteratively until all the outputs 158 of the CNN are summed appropriately. Once the final Distribution is calculated, it is simply a 159 matter of getting its logits which can be used to calculate the loss of the predictions. 160

161 DOLPHIN provides additional primitives to support more complex symbolic programs. Figure 3b, for instance, shows the computation graph for the PathFinder task (Tay et al.) 2021), which involves

162 $\mathbf{APPLY} \quad : \quad \mathbb{D}^K \times (S^K \to S) \to \mathbb{D}$ 163 $s \in S$ (objects) Filter : $\mathbb{D} \times (\hat{S} \to \mathbb{B}) \to \mathbb{D}$ Symbol :: 164 Tag :: $t \in T$ (tensors) APPLYIF : $\mathbb{D}^K \times (S^K \to S) \times (S^K \to \mathbb{B}) \to \mathbb{D}$ 165 Distribution :: $D \in \mathbb{D} = S \to T$ Union : $\mathbb{D} \times \mathbb{D} \to \mathbb{D}$ 166 GetProbs : $\mathbb{D} \to [0,1]^N$ 167

168

Figure 4: Formal definition of DOLPHIN's programming abstractions and primitives.

169 recursively building paths to identify if two points in a maze are connected. The union primitive is 170 used to support the recursive nature of this program. As with apply, DOLPHIN maps the symbolic 171 operations denoted within these primitives to probability computations. Given that Distribution objects associate symbols with the batched logits themselves, these computations are vectorized 172 and directly operate over PyTorch tensors. This deep integration of DOLPHIN into PyTorch allows 173 programmers to write symbolic programs as symbolic layers that interact with standard PyTorch 174 neural layers within a neurosymbolic model. DOLPHIN can thus leverage the hardware acceleration 175 supported by PyTorch to scale to large and complex programs. This contrasts with systems like Scal-176 lop (Li et al., 2023), where tensors are converted into Scallop-friendly tags transferred to a process 177 outside the Python environment with CPU-bound probability computations, restricting scalability. 178

179 180

181

182 183

184 185

187 188

189

190

191

192

193

194

197

3 THE DOLPHIN FRAMEWORK

3.1 DOLPHIN CORE DESIGN PRINCIPLES

We based DOLPHIN's framework design on the following core principles:

- Flexible programmability: The framework should allow developers to write neurosymbolic applications in Python with minimal effort, providing intuitive primitives that seamlessly integrate with Python's rich and expressive language features.
- End-to-end differentiability on GPUs: The framework should allow any neurosymbolic program to be end-to-end differentiable on GPUs irrespective of the task characteristics.
- Scalable: The framework should easily scale with greater problem and data complexity.
- **Tunable:** Similar to hyperparameters in deep learning, the framework should provide a simple interface for developers to choose provenances (and their configurations) or define new ones.

Together, these principles help address the challenges of scaling neurosymbolic frameworks. The flexible programmability and tunability allow us to write complex neurosymbolic programs, while 195 GPU differentiability and scalability work towards addressing data complexity. We show how we 196 realize these principles by describing the key components of DOLPHIN.

3.2 THE DOLPHIN SYNTAX 199

DOLPHIN provides a programming interface that developers can use to express symbolic programs 200 in a Pythonic manner. DOLPHIN maps each operation of the symbolic program to PyTorch 201 which enables end-to-end GPU-accelerated differentiable reasoning. Figure 4 presents DOLPHIN's 202 programming interface including the symbolic abstractions and operations over them. 203

204 3.2.1 Abstractions 205

206 The three main abstractions provided by DOLPHIN for expressing differentiable symbolic programs are shown on the left of Figure 4. Symbols S represent symbolic entities relevant to the program. 207 These entities can be any Pythonic object, such as hand-written digits in MNIST-SumN or coordi-208 nates of points in PathFinder. Tags T are tensors that represent their likelihoods. Typically, tags for 209 symbols are derived from the outputs of machine learning models, such as the probability distribu-210 tion over digits produced by the CNN classifier in MNIST-SumN. Finally, Distribution D represents 211 the likelihood of an input being classified as one of the pre-defined symbols. 212

Distributions serve as the fundamental datatype of a DOLPHIN program and act as its main interface 213 with a machine learning model. For instance, when the developer instantiates a Distribution object, 214 such as in the following code snippet from Figure 2a 215

Dres = Distribution(self.CNN(imgs[0]), digits)

the output of the CNN model is directly passed to the Distribution object, effectively acting as
 an input to the symbolic program. The Distribution object itself, as shown in Figure 2b, contains
 batches of *tags* extracted from the model outputs, and maintains the set of corresponding symbols.

To enable such a seamless integration between the PyTorch model and the symbolic program, Distributions are designed to operate directly over PyTorch tensors. This has two main advantages. First, it preserves the gradients of the model output throughout the symbolic program, enabling end-toend differentiability. Second, it allows DOLPHIN to perform operations over an entire batch of tags, leveraging the vectorized operations provided by PyTorch. DOLPHIN can thus operate efficiently on specialized hardware like GPUs, allowing the symbolic program to scale effectively.

3.2.2 OPERATIONS

Figure 4 shows the five operations supported by DOLPHIN that developers can use to manipulate Distributions and express complex symbolic programs. We now expand on these operations.

APPLY. This is the primary operation that developers can use to manipulate Distributions. It takes as inputs K Distributions, where $K \ge 1$, along with a function f of the same arity. This function defines operations over the symbols of K distributions. APPLY then computes the results of f over all possible combinations of arguments sourced from the symbols of the Distributions as well as their associated tags, and returns a new Distribution with these results and tags.

This operation occurs in two stages akin to the popular map-reduce pattern. In the *map* stage, APPLY computes the results of f over the symbols of the input Distributions and their associated tags:

$$R = \{ (f(s_1, s_2, \dots, s_k), (t_1 \otimes t_2 \otimes \dots \otimes t_k)) \mid D_i(s_i) = t_i, i = 1, \dots, k \}$$
(1)

Here, the tag of each result symbol $f(s_1, s_2, \ldots, s_k)$ is the conjunction \otimes of the tags (t_1, t_2, \ldots, t_k) of the input symbols it was derived from. While the tag computations are performed on the GPU, the function f is executed sequentially on the CPU for each combination of symbols. This is because function f can be any user-defined Python function, including complex control flows and operations like regex parsing, image processing, or Python's eval(). It may also be a many-to-one function and the tags shared by a resulting symbol must be aggregated to form the final tags of the output Distribution. We, therefore, *shuffle* the results from the map stage to compute a function M from each symbol to tags from R associated with it:

246 247

250

251

226

230

238

$$M = \lambda s \, \left\{ \, t \mid (s,t) \in R \, \right\} \tag{2}$$

We then proceed to the *reduce* stage, where we aggregate the tags of each symbol in M using disjunction \oplus to produce the final Distribution D_{res} :

$$D_{\rm res} = \lambda \ s \ . \bigoplus \left\{ \ t \mid t \in M(s) \right\}$$
(3)

Since the tags here are PyTorch tensors representing probabilities, the implementations of the conjunction and disjunction operations are dictated by the underlying provenance used by the program.
 A more detailed explanation of the provenances is provided in Section [3.4]

FILTER. The FILTER operation is used to filter out symbols from a Distribution based on some condition. It takes in a single Distribution, along with a user-defined function that returns a boolean value, which acts as the condition. This operation then returns a new Distribution that contains only the symbols that satisfy the condition, along with their tags.

APPLYIF. This operation is a conditional version of APPLY. It takes in K Distributions and functions f_{apply} and f_{cond} of the same arity. For each combination of symbols from the K Distributions, APPLYIF computes f_{apply} and its associated tags only if the condition f_{cond} is satisfied over that combination of symbols. The operation then returns a new Distribution with these results and tags.

UNION. The UNION operation is used to combine two Distributions. It takes in two Distributions and returns a new Distribution that contains the union of the symbols from the two input Distributions, along with their tags. Any symbols common to both input Distributions have their tags merged via a disjunction operation. UNION is especially useful when writing recursive programs in DOLPHIN that require combining the results of multiple recursive calls, as described in Appendix D.

GETPROBS. The GETPROBS operation extracts the probabilities from the tags of a Distribution. This is used mainly once the symbolic program has been executed to extract the final probabilities

Provenance	Domain	0	1	$t\oplus t'$	
DAMP	[0, 1]	0	1	$\operatorname{clamp}_0^1(t+t')$	$t \cdot t'$
DTKP-AM	$[0,1] \cup \{\infty,-\infty\}$	$\mathbf{\hat{0}}_{ij} = -\infty$	$\mathbf{\hat{1}}_{ij} = \begin{cases} \infty \\ -\infty \end{cases}$	$\begin{vmatrix} i = 1 \\ i > 1 \end{vmatrix} \operatorname{top}_k(\operatorname{cat}(t, t'))$	$\log_k([\min(t_i , t_j') \mid (t_i, t_j') \in t \times t$

Table 1: DOLPHIN provenances implemented in PyTorch.

of the symbols in the output Distribution. These probabilities can then be used to compute the loss function for training the neural model. The actual extraction of the probabilities from the tags depends on the specific provenance used in the program.

280 3.3 CONTROL FLOW AND RECURSION

Expressing control flow and recursion in a DOLPHIN program can be done in one of two ways.
The simplest way is to specify any control flow operations within the user-defined functions supplied to APPLY, APPLYIF, and FILTER, since these functions can contain arbitrary Python code.

Alternatively, one can specify control flow and re-285 cursion outside of these functions by specifying all 286 branches separately and merging their results using 287 UNION. Figure 5 shows an example of transitive 288 closure in DOLPHIN, where the compute_paths 289 function computes the transitive closure of the 290 graph by iteratively applying edges predicted by a 291 neural model to paths. The APPLYIF function ap-292 plies the edges to the paths if the end of the first 293 path is the same as the start of the second path. The UNION function merges the new paths with 294 the existing paths. The function compute_paths is 295 called recursively until a fixpoint is reached, specif-296 ically until no new paths can be added. 297



Figure 5: Transitive Closure in DOLPHIN.

298 299

275 276

277

278

279

281

3.4 DOLPHIN PROVENANCES

As discussed earlier, each symbol in a distribution is associated with a batch of one or more tags. The DOLPHIN primitives define how to manipulate certain tags. For instance, Equations (1) and (3) specify the tags to be conjuncted or disjuncted. We now define the semantics of these operations.

The goal of such operations is to approximate the probabilities of the symbols in the output distribution as accurately as possible. This is achieved by using a mathematical framework called *provenance semirings* (Green et al., 2007). Provenance semirings provide generalized algebraic structure to propagate tags when computing over tagged data. In the case of DOLPHIN distributions, we can view the tags as representing the probabilities, and the data as the distribution's symbols.

Designing and implementing provenances can be challenging since they must be accurate enough to capture the semantics of the symbolic program, while at the same time being coarse enough to maintain computational feasibility. Furthermore, the provenances must be differentiable to enable end-to-end training for neurosymbolic tasks. While neurosymbolic frameworks like Scallop (Li et al.) 2023) implement differentiable provenances, they are not designed to leverage hardware accelerations or batched optimizations due to the CPU-bound nature of their implementations. We thus design vectorized provenances in DOLPHIN that are differentiable and enable GPU computations.

We simplify the definition of provenances from Scallop as a 5-tuple: $(T, 0, 1, \otimes, \oplus)$. Here, T is the tag space, $\otimes : T \times T \to T$ is the conjunction operator with identity 0, and $\oplus : T \times T \to T$ is the disjunction operator with identity 1. We then implement two differentiable provenances in DOLPHIN: Differentiable Add-Mult Probabilities (DAMP) and Differentiable Top-K Proofs (DTKP). Table 1 summarizes the operations of these provenances.

Differentiable Add-Mult Probabilities. Differentiable Add-Mult Probabilities (DAMP) is a popular technique that uses the probability space as its tag space: T = [0, 1]. Its conjunction operation \otimes is defined as the product of probabilities, clamped at 1, and its disjunction operation \oplus is defined as the sum of probabilities. The main assumption underlying the DAMP operations is that the input Distributions are mutually exclusive and independent. This assumption allows DAMP to compute probabilities extremely efficiently, as the operations are simple and can be easily vectorized.

Differentiable Top-k Proofs. Differentiable Top-k Proofs (DTKP) (Huang et al., 2021) was pro-327 posed to overcome the shortcomings of DAMP. This provenance tracks a set of up to k proofs for 328 each symbol. Each proof, in turn, denotes the set of input symbols necessary to derive the output symbol. These proofs are then used to compute the probabilities of the output symbols. In Scal-330 lop, DTKP tags are converted into probabilities via differentiable weighted model counting (WMC). 331 This form of DTKP, which we call DTKP-WMC, is computationally hard and is by nature difficult 332 to vectorize due to the varying sizes of proof sets and the WMC procedure. We hence design a 333 vectorized approximation of DTKP-WMC, called DTKP-AM (DTKP with Add-Mult), that can be 334 efficiently computed on GPUs.

We first define the structure of tags in DTKP-AM in a manner that conforms to the constraints of PyTorch tensors. Each tag t for a symbol s is a 2-dimensional tensor of shape (k, |I|), where k is the maximum number of proofs to be retained and I is an ordered list of all *input symbols* (symbols that are present in the input Distributions). Each row t_i of t corresponds to one of the tag's k proofs. Each element t_{ij} thus represents the probability of the *j*th input symbol in the *i*th proof:

$$t_{ij} = \begin{cases} p_j & \text{if the } j \text{th symbol is present in the } i \text{th proof} \\ \hat{\mathbf{0}}_{ij} & \text{otherwise} \end{cases}$$

where p_j is the probability of the *j*th input symbol. The probability of each proof is then computed by taking the product of the normal:

$$\Pr(t_i) = \prod_j \operatorname{norm}(t_{ij}) \quad \text{where} \quad \operatorname{norm}(t_{ij}) = \begin{cases} 1 & t_{ij} = +\infty \\ 0 & t_{ij} = -\infty \\ t_{ij} & \text{otherwise} \end{cases}$$

We next define the operations of DTKP-AM in Table []. The \oplus operation is defined as the union of two tag tensors t and t' while \otimes is defined as the element-wise minimum of the normalized elements of all possible combinations of proofs in t and t'. In each case, the top_k operation retains only upto k proofs with the highest probabilities. These definitions thus allow us to take advantage of the benefits of the DTKP provenance while enabling efficient computation on GPUs. To calculate the probability of the entire tag, DTKP-AM adds the probabilities of the individual proofs and clamps it at 1. We provide a detailed discussion of DTKP-AM in Appendix C

354 355 356

340

341 342

343

3.5 BUILDING THE DOLPHIN PROGRAM

The programmer specifies the neurosymbolic task using a Python program P that uses DOLPHIN's programming interface to connect the neural components (e.g., neural networks) with the symbolic components and operations. We call P the *symbolic program*. Because P is a Python program and DOLPHIN interfaces with PyTorch, DOLPHIN supports any PyTorch-based neural network(s), most Python language features, and custom user-defined functions. This feature enables greater flexibility and expressiveness in neurosymbolic programs than existing frameworks.

In addition to P, the programmer provides one or more neural networks M_1, \ldots, M_k , and a dataset D. Given these inputs, DOLPHIN extracts the computation graph that encodes how the neural network outputs are transformed using symbolic operations to produce a final result D_{res} . All computations in DOLPHIN are expressed using distribution objects D_i . Each DOLPHIN primitive (e.g., APPLY) takes one or more distribution objects as inputs and applies a transformation to produce another distribution object. Because each distribution object D_i only contains vectors of tags (or probabilities), the entire computation graph (including the neural network(s)) can be ported to a GPU for faster execution. During training, DOLPHIN optimizes over the standard objective function:

$$\phi(\theta) = \min_{\theta} \sum_{(x,y)\in\mathcal{D}} \mathcal{L}(P(M_{\theta}(x)), y)$$
(4)

Here \mathcal{L} is the loss function, such as binary cross entropy. While DOLPHIN allows P to take multiple neural networks as inputs, we show only one neural network model M here for simplicity.

370

371 372

4 Experiments

377 We evaluate DOLPHIN on a set of 13 benchmarks of varying complexity and scale across 5 neurosymbolic tasks. Our evaluation addresses the following research questions:

- **RQ1: Scalability.** How does DOLPHIN scale to complex problems and large datasets?
 - **RQ2:** Accuracy. Do models written in DOLPHIN converge to SOTA accuracies?
 - RQ3: Provenance Comparisons. Which provenances are most effective for each benchmark?

382383 4.1 BENCHMARKS

380

381

We evaluate DOLPHIN on the following benchmarks. We give additional context and information about the experiment setup for each branch in Appendix A

387 **MNIST-SumN.** The MNIST-SumN (or briefly, SumN) task from (De Smet et al., 2024) takes as 388 inputs N handwritten digits from the MNIST dataset and returns their sum. We consider three 389 versions of this task: small (N = 5), medium (N = 10), and large (N = 15).

Hand-Written Formula (HWF). The HWF task from Li et al. (2020) takes as input a set of images of handwritten digits and arithmetic operators representing a formula. The task is to evaluate the formula and return the result. We consider three versions of HWF: small (formulas of length up to 7), medium (formulas of length up to 15), and large (formulas of length up to 19).

PathFinder. PathFinder (or Path) from [Tay et al.] (2021) tests the ability of an agent to reason over long-range dependencies within an input image. The image consisting of two dots and a sequence of curved and dashed lines. The task is to identify whether the two dots are connected via the lines.
We consider three versions of this task based on the image size in pixels: small (32 x 32), medium (128 x 128), and large (256 x 256).

CLUTRR. In this task from Sinha et al. (2019), the input is a passage of text containing some information about several individuals and some of their relationships. The task is then to infer the relationship between two given individuals, which is not explicitly provided in the input. We consider two versions of this task, where the training data contains relation chains of lengths up to 3 (small) or 4 (medium).

Mugen. In this task from Hayes et al. (2022), the input is a video of gameplay footage that is 3.2 seconds long and a natural language passage captioning the video. The goal is to measure how aligned the text is with the video. This task has two variants: Mugen-TVR, where the model must retrieve the video that best aligns with the text, and Mugen-VTR, where the model must retrieve the text that best aligns with the video. This benchmark tests the ability of the model to reason over multimodal data. We consider two versions of this task: small, with 1000 training samples, and medium, with 5000 training samples.

411 412

413

4.2 EXPERIMENTAL SETUP AND BASELINES

Setup. All experiments, except for the CLUTRR benchmark, were run on machines with two
20-core Intel Xeon Gold 6248 CPUs, four NVIDIA GeForce RTX 2080 Ti GPUs, and 768 GB
RAM. Since the CLUTRR benchmark requires more GPU memory, it was run on a machine with 8
NVIDIA A100 40GB GPUs instead. We ran each tool thrice until convergence or until a timeout of
10 hours was reached and report the average best accuracy and training time. The code of DOLPHIN
and the benchmarks are provided in the supplementary material.

Baselines. We select Scallop (Li et al., 2023), a contemporary state-of-the-art neurosymbolic framework supporting differentiable programming optimized to run on the CPU and use multiple cores to
parallelize its computations. We also choose two sampling-based gradient approximation methods,
ISED (Solko-Breslin et al., 2024) and IndeCateR+ (De Smet et al., 2024). We compare DOLPHIN
against Scallop on all benchmarks, and against ISED and IndeCateR on MNIST-SumN and HWF.

425 426 4.3 RQ1: SCALABILITY

Table 2 presents the total training times (T_{total}) in seconds for DOLPHIN and baselines on all benchmarks, along with the time per epoch (T_{epoch}) and scaling factor α . The scaling factor is the ratio of the per epoch times of the baselines to DOLPHIN. We observe that in almost all cases, DOLPHIN shows significant improvements in training times, training up to about 300x faster than Scallop, 44x faster than ISED, and 3x faster than IndeCateR+. On average, DOLPHIN reports a speedup of 21x times over all baselines. Further, DOLPHIN converges on all benchmarks, while the baselines time

Task	Dol	PHIN		Scallop		ISED			IndeCateR+		
	T _{total}	Tepoch	T _{total}	T _{epoch}	α	T _{total}	Tepoch	α	T _{total}	T _{epoch}	
SumN (S)	78.33	15.67	923.78	184.76	11.80	299.63	59.93	3.82	416.78	59.54	3
SumN (M)	144.92	14.49	3.41e3	341.57	23.57	2.16e3	216.54	14.94	385.65	32.14	2.3
SumN (L)	220.47	14.70	7.41e3	493.87	33.60	9.8e3	653.39	44.45	548.28	23.84	1.
HWF (S)	3.17e3	158.79	9.99e3	499.57	3.15	1.58e4	790.04	4.97	1.35e4	540.26	3
HWF (M)	1.46e4	731.67	TO	1.49e4	20.41	TO	6.83e3	9.34	2.51e4	2512.03	3.
HWF (L)	2.42e4	1.21e3	ТО	3.92e5	323.21	ТО	1.05e4	8.61	ТО	4091.46	3.
Path (S)	1.08e4	1.08e3	2.2e4	2.2e3	2.03	N.A.					
Path (M)	1.79e4	1.79e3	TO	4.17e3	2.32						
Path (L)	1.94e4	1.94e3	ТО	1.12e4	5.81						
CLUTRR (S)	1.54e3	154.85	4.29e3	429.97	2.77	N.A.					
CLUTRR (M)	2.91e3	291.36	7.83e3	783.11	2.69						
Mugen (S)	3.62e3	180.80	1.34e4	133.68	0.74						
Mugen (M)	1.78e3	890.34	TO	634.86	0.71			IN.	А.		

Table 2: Comparison of training times taken by each baseline. The Timeout (TO) is set at 10 hours. α is the scaling factor, which is the ratio of the per epoch training times of the baselines and DOLPHIN.



Figure 6: Accuracy of DOLPHIN and baselines across all benchmarks.

out on most of the Medium and Large versions of the benchmarks. These results thus demonstrate that, unlike existing tools, DOLPHIN can scale to complex problems and large datasets.

DOLPHIN trains slightly slower than Scallop on both versions of the Mugen benchmark. This is because the DOLPHIN program written for Mugen uses Python objects and operations that are not fully batchable across samples. In contrast, the Scallop program, which is written in a compiled and optimized language, runs around 1.3x faster than DOLPHIN on average per iteration. However, DOLPHIN requires only 20 epochs to converge whereas Scallop requires almost 1000 epochs (Li et al., 2023). As a result, DOLPHIN's total training time is still significantly lower than Scallop's (3.7x for Mugen(S)). We show the accuracy curves for Mugen in Appendix A.6

4.4 RQ2: ACCURACY

Figure 6 presents the accuracy of DOLPHIN and the baselines on the different benchmarks. DOL-PHIN accuracies are marked in blue. In all cases, for DOLPHIN, we report the accuracies of the best-performing provenance. We use the DAMP provenance for MNIST, CLUTRR, and Mugen benchmarks, and the DTKP-AM provenance for the HWF and PathFinder benchmarks.

We observe that in all cases, DOLPHIN achieves state-of-the-art accuracy among neurosymbolic frameworks, except in CLUTRR, where DOLPHIN's accuracy is slightly lower than Scallop's. While DeepProbLog (Manhaeve et al., 2018) reports near-perfect accuracies for CLUTRR, they use nega-tive mining techniques to provide additional labels at train time. Scallop and DOLPHIN, on the other hand, stick to a traditional semi-supervised multiclass classification approach. For most Medium and Large versions of the benchmarks, DOLPHIN achieves better accuracy, whereas the baselines either report lower accuracy due to the complexity of the benchmark (e.g., HWF) or fail to converge within the time limit (e.g., Scallop on PathFinder-Large). Most importantly, these results show that DOLPHIN's scalability improvements do not come at the cost of accuracy.

4.5 RQ3: PROVENANCE COMPARISONS

We perform ablation studies to compare the effectiveness of the DAMP and DTKP-AM provenances for each benchmark. We share the graphs in Figure 9 (Appendix B). In all cases, training with the DAMP provenance takes around 132.96 seconds per epoch less than with DTKP-AM on average.

However, the effectiveness of each provenance varies from benchmark to benchmark. For all variations of CLUTRR, Path, and Mugen, both provenances achieve comparable accuracies, with DTKP-AM usually having a slight edge. In the MNIST-SumN benchmark, the DAMP provenance is more effective than the DTKP-AM provenance by 72.08 %pts on average, since the top-k proofs cannot capture all the possible ways in which the sum of the digits can be computed.

In contrast, for HWF, the DTKP-AM provenance is more effective than DAMP by an average of
42.18 %pts. Each step of the HWF program, shown in Appendix G, involves both a concatenation
operation and a partial parsing operation before the final expression is evaluated to produce a result.
As such, it is difficult for the tags in DAMP to capture the semantics of the symbolic program. In
the case of DTKP-AM, each tag is a collection of proofs over input symbols corresponding to logits
derived from the neural model. Therefore, any calculated gradients can be directly backpropagated
to the logits that most influenced the output, making this a more effective provenance for this task.

498 499

5 RELATED WORK

500 Neurosymbolic programming frameworks. Frameworks like Scallop (Li et al., 2023), Deep-ProbLog (Manhaeve et al., 2018), and ISED (Solko-Breslin et al., 2024) provide a simple interface 501 for neurosymbolic programming. There are also domain-specific tools like NeurASP Yang et al. 502 (2021) for answer set programming and NeuralLog Chen et al. (2021) for phrase alignment in NLP. 503 While these frameworks provide intuitive abstractions, they are bottle-necked due to expensive data 504 transfers between symbolic computations done on CPU versus neural computations that execute on 505 GPU, making neurosymbolic learning hard to scale. In contrast, DOLPHIN provides a deeper inte-506 gration of the two worlds by building a Python-based API on top of PyTorch, which scales better. 507

Scaling techniques. Several optimization techniques have been proposed to improve the scalability 508 of differentiable reasoning algorithms. Some techniques aim to scale reasoning algorithms by com-509 piling the symbolic program into computation graphs that can be run on GPUs. LYRICS (Marra 510 et al., 2019), Logic Tensor Networks (Badreddine et al., 2022), and Tensorlog (Cohen et al., 2020) 511 are examples of such techniques. However, these methods focus on first order logic programs and 512 provide limited support for user-defined Pythonic functions, essential for building complex neu-513 rosymbolic programs, Greedy NTP (Minervini et al., 2020a) reduces the computation cost of 514 NTP (Rocktäschel & Riedel, 2017) by tracking only a subset of proof states using nearest neigh-515 bor search. Likewise, the conditional theorem prover (Minervini et al., 2020b) employs a machine 516 learning-based proof selection technique. However, unlike DOLPHIN, these methods are point solu-517 tions that do not fundamentally address the scalability challenge for neurosymbolic learning.

518 **Specialized neurosymbolic solutions.** There are many specialized solutions for various neurosym-519 bolic tasks. For instance, NGS (Li et al., 2020) uses a hand-coded syntax to specify the structure 520 of mathematical expressions for HWF. More general solutions, such as NS-CL (Mao et al., 2019) 521 includes a framework for visual question answering that learns symbolic representations for text and 522 images. NeRd (Chen et al., 2021) transforms questions in natural language into executable programs 523 based on symbolic information extracted from text. Orvieto et al. (2023) proposes a recurrent neural network architecture that achieves 95% accuracy on Path (S) and 94% on Path (M). In contrast, 524 Dolphin is a general programming framework that tries to scale diverse neurosymbolic programs. 525

526 527

6 CONCLUSION AND LIMITATIONS

We proposed DOLPHIN, a programmable framework for scaling neurosymbolic learning. DOLPHIN provides abstractions for writing symbolic programs along with pluggable vectorized provenances to compute symbolic gradients. This allows users to write differentiable symbolic programs in Python within PyTorch pipelines that can scale to complex programs and large datasets. We show that DOLPHIN scales significantly better than existing neurosymbolic frameworks while achieving state-of-the-art performance on a variety of tasks.

A limitation of DOLPHIN is that it needs the user to write programs in a batched manner. While this is a common pattern within deep learning, it may be restrictive to users new to batched programming. Also, while DOLPHIN works well with most models, the representation needed by generative models (e.g., Causal LLMs like Llama) has not been investigated. A third limitation is that Dolphin lacks support for non-deterministic symbolic programs. We leave this to future work.

539

¹Refer to Appendix E for a detailed comparison of DOLPHIN with LTN.

540 REFERENCES 541

549

577

583

584

585

591

542	Samy Badreddine, Artur d'Avila Garcez, Luciano Serafini, and Michael Spranger. Logic tensor net-
543	works. Artificial Intelligence, 303:103649, 2022. ISSN 0004-3702. doi: https://doi.org/10.1016/
544	j.artint.2021.103649. URL https://www.sciencedirect.com/science/article/
545	pii/S0004370221002009.

- 546 Swarat Chaudhuri, Kevin Ellis, Oleksandr Polozov, Rishabh Singh, Armando Solar-Lezama, Yisong 547 Yue, et al. Neurosymbolic programming. Foundations and Trends® in Programming Languages, 548 7(3):158-243, 2021.
- Zeming Chen, Qiyue Gao, and Lawrence S. Moss. NeuralLog: Natural language inference with joint 550 neural and logical reasoning. In Lun-Wei Ku, Vivi Nastase, and Ivan Vulić (eds.), Proceedings 551 of *SEM 2021: The Tenth Joint Conference on Lexical and Computational Semantics, pp. 78-552 88, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021. 553 starsem-1.7. URL https://aclanthology.org/2021.starsem-1.7. 554
- 555 William W. Cohen, Fan Yang, and Kathryn Mazaitis. Tensorlog: A probabilistic database im-556 plemented using deep-learning infrastructure. J. Artif. Intell. Res., 67:285-325, 2020. URL https://api.semanticscholar.org/CorpusID:211263674 558
- Meihua Dang, Pasha Khosravi, Yitao Liang, Antonio Vergari, and Guy Van den Broeck. Juice: A 559 julia package for logic and probabilistic circuits. In AAAI Conference on Artificial Intelligence, 560 2021. URL https://api.semanticscholar.org/CorpusID:235363700. 561
- 562 Adnan Darwiche. An advance on variable elimination with applications to tensor-based computa-563 tion. In ECAI 2020, pp. 2559–2568. IOS Press, 2020.
- 565 Lennert De Smet, Emanuele Sansone, and Pedro Zuidberg Dos Martires. Differentiable sampling of categorical distributions using the catlog-derivative trick. Advances in Neural Information 566 567 Processing Systems, 36, 2024.
- 568 Todd J Green, Grigoris Karvounarakis, and Val Tannen. Provenance semirings. In Proceedings of 569 the twenty-sixth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems, 570 pp. 31-40, 2007. 571
- 572 Thomas Hayes, Songyang Zhang, Xi Yin, Guan Pang, Sasha Sheng, Harry Yang, Songwei Ge, 573 Qiyuan Hu, and Devi Parikh. Mugen: A playground for video-audio-text multimodal understand-574 ing and generation. In European Conference on Computer Vision, pp. 431–449. Springer, 2022. 575
- Jiani Huang, Ziyang Li, Binghong Chen, Karan Samel, Mayur Naik, Le Song, and Xujie Si. Scallop: 576 From probabilistic deductive databases to scalable differentiable reasoning. Advances in Neural Information Processing Systems, 34:25134–25145, 2021. 578
- 579 Subbarao Kambhampati, Karthik Valmeekam, Lin Guan, Mudit Verma, Kaya Stechly, Siddhant 580 Bhambri, Lucas Paul Saldyt, and Anil B Murthy. Position: Llms can't plan, but can help planning 581 in llm-modulo frameworks. In Forty-first International Conference on Machine Learning. 582
 - Qing Li, Siyuan Huang, Yining Hong, Yixin Chen, Ying Nian Wu, and Song-Chun Zhu. Closed loop neural-symbolic learning via integrating neural perception, grammar parsing, and symbolic reasoning. In International Conference on Machine Learning, pp. 5884–5894. PMLR, 2020.
- 586 Ziyang Li, Jiani Huang, and Mayur Naik. Scallop: A language for neurosymbolic programming. 587 Proceedings of the ACM on Programming Languages, 7(PLDI):1463–1487, 2023. 588
- 589 Yinhan Liu. Roberta: A robustly optimized bert pretraining approach. arXiv preprint 590 arXiv:1907.11692, 2019.
- Robin Manhaeve, Sebastijan Dumancic, Angelika Kimmig, Thomas Demeester, and Luc De Raedt. 592 Deepproblog: Neural probabilistic logic programming. Advances in neural information processing systems, 31, 2018.

594 595 596 597 598	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 <i>International Conference on Learning Representations</i> , 2019. URL https://openreview.learning.com net/forum?id=rJgMlhRctm.
599 600 601 602 603 604	Giuseppe Marra, Francesco Giannini, Michelangelo Diligenti, and Marco Gori. Lyrics: A gen- eral interface layer to integrate logic inference and deep learning. In <i>Machine Learning and</i> <i>Knowledge Discovery in Databases: European Conference, ECML PKDD 2019, Würzburg,</i> <i>Germany, September 16–20, 2019, Proceedings, Part II</i> , pp. 283–298, Berlin, Heidelberg, 2019. Springer-Verlag. ISBN 978-3-030-46146-1. doi: 10.1007/978-3-030-46147-8_17. URL https://doi.org/10.1007/978-3-030-46147-8_17.
605 606 607	Pasquale Minervini, Matko Bošnjak, Tim Rocktäschel, Sebastian Riedel, and Edward Grefenstette. Differentiable reasoning on large knowledge bases and natural language. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 34, pp. 5182–5190, 2020a.
608 609 610	Pasquale Minervini, Sebastian Riedel, Pontus Stenetorp, Edward Grefenstette, and Tim Rocktäschel. Learning reasoning strategies in end-to-end differentiable proving. In <i>International Conference</i> <i>on Machine Learning</i> , pp. 6938–6949. PMLR, 2020b.
612 613 614 615	Aaditya Naik, Adam Stein, Yinjun Wu, Mayur Naik, and Eric Wong. Torchql: A program- ming framework for integrity constraints in machine learning. <i>Proc. ACM Program. Lang.</i> , 8(OOPSLA1), April 2024. doi: 10.1145/3649841. URL https://doi.org/10.1145/ 3649841.
616 617 618	Antonio Orvieto, Samuel L Smith, Albert Gu, Anushan Fernando, Caglar Gulcehre, Razvan Pas- canu, and Soham De. Resurrecting recurrent neural networks for long sequences. In <i>Proceedings</i> of the 40th International Conference on Machine Learning, ICML'23. JMLR.org, 2023.
619 620 621	Tim Rocktäschel and Sebastian Riedel. End-to-end differentiable proving. Advances in neural information processing systems, 30, 2017.
622 623	V Sanh. Distilbert, a distilled version of bert: Smaller, faster, cheaper and lighter. <i>arXiv preprint arXiv:1910.01108</i> , 2019.
624 625 626	Scallop Language Group. Scallop and neuro-symbolic programming: Tags, instrumentation, and provenance. Eleventh Summer School on Formal Techniques, 2022. URL https://www.scallop-lang.org/ssft22/lectures/lecture-2.pdf .
628 629 630 631 632 633 634	Koustuv Sinha, Shagun Sodhani, Jin Dong, Joelle Pineau, and William L. Hamilton. CLUTRR: A diagnostic benchmark for inductive reasoning from text. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pp. 4506–4515, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1458. URL https://aclanthology.org/D19-1458 .
635 636	Alaia Solko-Breslin, Seewon Choi, Ziyang Li, Neelay Velingker, Rajeev Alur, Mayur Naik, and Eric Wong. Data-efficient learning with neural programs. <i>arXiv preprint arXiv:2406.06246</i> , 2024.
637 638 639 640 641	Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena : A benchmark for efficient transformers. In <i>International Conference on Learning Representations</i> , 2021. URL https://openreview.net/forum?id=qVyeW-grC2k.
642 643 644 645 646 647	Zhun Yang, Adam Ishay, and Joohyung Lee. Neurasp: embracing neural networks into answer set programming. In Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI'20, 2021. ISBN 9780999241165.