

DISCRETE NEURAL ALGORITHMIC REASONING

Anonymous authors

Paper under double-blind review

ABSTRACT

Neural algorithmic reasoning aims to capture computations with neural networks via learning the models to imitate the execution of classic algorithms. While common architectures are expressive enough to contain the correct model in the weights space, current neural reasoners are struggling to generalize well on out-of-distribution data. On the other hand, classic computations are not affected by distributional shifts as they can be described as transitions between discrete computational states. In this work, we propose to force neural reasoners to maintain the execution trajectory as a combination of finite predefined states. To achieve that, we separate discrete and continuous data flows and describe the interaction between them. Trained with supervision on the algorithm’s state transitions, such models are able to perfectly align with the original algorithm. To show this, we evaluate our approach on multiple algorithmic problems and get perfect test scores both in single-task and multitask setups. Moreover, the proposed architectural choice allows us to prove the correctness of the learned algorithms for any test data.

1 INTRODUCTION

Learning to capture algorithmic dependencies in data and to perform algorithmic-like computations with neural networks are core problems in machine learning, studied for a long time using various approaches (Roni Kharon, 1994; Graves et al., 2014; Zaremba & Sutskever, 2014; Reed & De Freitas, 2015; Kaiser & Sutskever, 2015; Veličković et al., 2020b).

Neural algorithmic reasoning (Veličković & Blundell, 2021) is a research area focusing on building models capable of executing classic algorithms. Relying on strong theoretical guarantees of algorithms to work correctly on any input of any size and distribution, this setting provides unlimited challenges for out-of-distribution generalization of neural networks. Prior work explored this setup using the CLRS-30 benchmark (Veličković et al., 2022), which covers classic algorithms from the Introduction to Algorithms textbook (Cormen et al., 2009) and uses graphs as a universal tool to encode data of various types. Importantly, CLRS-30 also provides the decomposition of classic algorithms into subroutines and simple transitions between consecutive execution steps, called hints, which can be used during training in various forms.

The core idea of the CLRS-30 benchmark is to understand how neural reasoners generalize well beyond the training distribution, namely on larger graphs. Classic algorithms possess strong generalization due to the guarantee that correct execution steps never encounter ‘out-of-distribution’ states, as all state transitions are predefined by the algorithm. In contrast, when encountering inputs from distributions that significantly differ from the train data, neural networks are usually not capable of robustly maintaining internal calculations in the desired domain. Consequently, due to the complex and diverse nature of all possible data that neural reasoners can be tested on, the generalization performance of such models can vary depending on particular test distribution (Mahdavi et al., 2023).

Given that, it is becoming important to interpret internal computations of neural reasoners to find errors or to prove the correctness of the learned algorithms (Georgiev et al., 2021).

Interpretation methods have been actively developing recently due to various real-world applications of neural networks and the need to debug and maintain systems based on them. Especially, the Transformer architecture (Vaswani et al., 2017) demonstrates state-of-the-art performance in natural language processing and other modalities, representing a field for the development of interpretability

054 methods (Elhage et al., 2021; Weiss et al., 2021; Zhou et al., 2024; Lindner et al., 2024). Based on
055 active research on a computational model behind the transformer architecture, recent works propose
056 a way to learn models that are fully interpretable by design (Friedman et al., 2023).

057 We found the ability to design models that are interpretable in a simple and formalized way to be
058 crucial for neural algorithmic reasoning as it is naturally related to the goal of learning to perform
059 computations with neural networks.

060
061 In this paper, we propose to force neural reasoners to follow the execution trajectory as a combina-
062 tion of finite predefined states, which is important for both generalization ability and interpretability
063 of neural reasoners. To achieve that, we start with an attention-based neural network and describe
064 three building blocks to enhance its generalization abilities: feature discretization, hard attention
065 and separating discrete and continuous data flows. In short, all mentioned blocks are connected:

- 066 • State discretization does not allow the model to use complex and redundant dependencies
067 in data;
- 068 • Hard attention is needed to ensure that attention weights will not be annealed for larger
069 graphs. Also, hard attention limits the set of possible messages that each node can receive;
- 070 • Separating discrete and continuous flows is needed to ensure that state discretization does
071 not lose information about continuous data.

072
073 Then, we build fully discrete neural reasoners for different algorithmic tasks and demonstrate their
074 ability to perfectly mimic ground-truth algorithm execution. As a result, we achieve perfect test
075 scores on the multiple algorithmic tasks with guarantees of correctness on any test data. Moreover,
076 we demonstrate that a single network is capable of executing all covered algorithms in a multitask
077 manner with perfect generalization too.

078 In summary, we consider the proposed blocks as a crucial component for robust and interpretable
079 neural reasoners and demonstrate that trained with hint supervision, discretized models perfectly
080 capture the dynamic of the underlying algorithms and do not suffer from distributional shifts.

081 082 2 BACKGROUND

083 084 2.1 ALGORITHMIC REASONING

085
086 Performing algorithmic-like computations usually requires the execution of sequential steps and the
087 number of such steps depends on the input size. To imitate such computations, neural networks are
088 expected to be based on some form of recurrent unit, which can be applied to a particular problem
089 instance several times (Kaiser & Sutskever, 2015; Zaremba & Sutskever, 2014; Vinyals et al., 2015;
090 Veličković et al., 2020b).

091 The CLRS Algorithmic Reasoning Benchmark (CLRS-30) (Veličković et al., 2022) defines a general
092 paradigm of algorithmic modeling based on Graph Neural Networks (GNNs), as graphs can natu-
093 rally represent different input types and manipulations over such inputs. Also, GNNs are proven to
094 be well-suited for neural execution (Xu et al., 2020; Dudzik & Veličković, 2022).

095 The CLRS-30 benchmark covers different algorithms over various domains (arrays, strings, graphs)
096 and formulates them as algorithms over graphs. Also, CLRS-30 proposes to utilize the decom-
097 position of the algorithmic trajectory execution into simple logical steps, called hints. Using this
098 decomposition is expected to better align the model to desired computations and prevent it from
099 utilizing hidden non-generalizable dependencies of a particular train set. Prior work demonstrates a
100 wide variety of additional inductive biases for models towards generalizing computations, including
101 different forms of hint usage (Veličković et al., 2022; Bevilacqua et al., 2023), biases from standard
102 data structures (Jürß et al., 2024; Jain et al., 2023), knowledge transfer and multitasking (Xhonneux
103 et al., 2021; Ibarz et al., 2022; Numeroso et al., 2023), etc. Also, recent studies demonstrate sev-
104 eral benefits of learning neural reasoners end-to-end without any hints at all (Mahdavi et al., 2023;
105 Rodionov & Prokhorenkova, 2023).

106 The recently proposed SALSA-CLRS benchmark (Minder et al., 2023) enables a more thorough
107 OOD evaluation compared to CLRS-30 with increased test sizes (up to 100-fold train-to-test scaling,
compared to 4-fold for CLRS-30) and diverse test distributions. Despite significant gains in the

108 performance of neural reasoners in recent work, current models still struggle to generalize to out-
 109 of-distribution (OOD) test data (Mahdavi et al., 2023; Georgiev et al., 2023; Minder et al., 2023).
 110 While de Luca & Fountoulakis (2024) prove by construction the ability of the transformer-based
 111 neural reasoners to perfectly simulate graph algorithms (with minor limitations occurring from the
 112 finite precision), it is still unclear if generalizable and interpretable models can be obtained via
 113 learning. Importantly, the issues of OOD generalization are induced not only by the challenges
 114 of capturing the algorithmic dependencies in the data but also by the need to carefully operate
 115 with continuous inputs. For example, investigating the simplest scenario of learning to emulate
 116 the addition of real numbers, Klindt (2023) demonstrates the failure of some models to exactly
 117 imitate the desired computations due to the nature of gradient-based optimization. This limitation
 118 can significantly affect the performance of neural reasoners on adversarial examples and larger input
 119 instances when small errors can be accumulated.

120 2.2 TRANSFORMER INTERPRETABILITY AND COMPUTATION MODEL

122 Transformer (Vaswani et al., 2017) is a neural network architecture for processing sequential data.
 123 The input to the transformer is a sequence of tokens from a discrete vocabulary. The input layer maps
 124 each token to a high-dimensional embedding and aggregates it with the positional encoding. The
 125 key components of each layer are attention blocks and MLP with residual connections. Providing a
 126 detailed description of mechanisms learned by transformer models (Elhage et al., 2021) is of great
 127 interest due to their widespread applications.

128 RASP (Weiss et al., 2021) is a programming language proposed as a high-level formalization of the
 129 computational model behind transformers. The main primitives of RASP are elementwise sequence
 130 functions, *select* and *aggregate* operations, which conceptually relate to computations performed
 131 by different blocks of the model. Later, Lindner et al. (2024) presented Tracr, a compiler for con-
 132 verting RASP programs to the weights of the transformer model, which can be useful for evaluating
 133 interpretability methods.

134 While RASP might have limited expressibility, it supports arbitrary complex continuous functions
 135 which in theory can be represented by transformer architecture, but are difficult to learn. Also, RASP
 136 is designed to formalize computations over sequences of fixed length. Motivated by that, Zhou et al.
 137 (2024) proposed RASP-L, a restricted version of RASP, which aims to formalize the computations
 138 that are easy to learn with transformers in a size-generalized way. The authors also conjecture that
 139 the length-generalization of transformers on algorithmic problems is related to the ‘simplicity’ of
 140 solving these problems in RASP-L language.

141 Another recent work (Friedman et al., 2023) describes Transformer Programs: constrained trans-
 142 formers that can be trained using gradient-based optimization and then automatically converted into
 143 a discrete, human-readable program. Built on RASP, transformer programs are not designed to be
 144 size-invariant.

146 3 DISCRETE NEURAL ALGORITHMIC REASONING

147 3.1 ENCODE-PROCESS-DECODE PARADIGM

149 Our work follows the encode-process-decode paradigm (Hamrick et al., 2018), which is usually
 150 employed for step-by-step neural execution.

151 All input data is represented as a graph G with an adjacency matrix A and node and edge features
 152 that are first mapped with a simple linear encoder to high-dimensional vectors of size h . Let us
 153 denote node features at a time step t ($1 \leq t \leq T$) as $X^t = (x_1^t, \dots, x_n^t)$ and edge features as
 154 $E_t = (e_1^t, \dots, e_m^t)$. Then, the processor, usually a single-layer GNN, recurrently updates these
 155 features, producing node and edge features for the next step:

$$156 X^{t+1}, E^{t+1} = \text{Processor}(X^t, E^t, A).$$

157
 158 The processor network can operate on the original graph defined by the task (for graph problems) or
 159 on the fully connected graph. For the latter option, the information about the original graph can be
 160 encoded into the edge features.
 161

The number of processor steps T can be defined automatically by the processor or externally (e.g., as the number of steps of the original algorithm). After the last step, the node and edge features are mapped with another linear layer, called the decoder, to the output predictions of the model.

If the model is trained with hint supervision, the changes of node and edge features at each step are expected to be related to the original algorithm execution. In this sense, the processor network is aimed to mimic the algorithm’s execution in the latent space.

3.2 DISCRETE NEURAL ALGORITHMIC REASONERS

In this section, we describe the constraints for the processor that allow us to achieve a fully interpretable neural reasoner. We start with Transformer Convolution (Shi et al., 2020) with a single attention head.

As mentioned above, at each computation step t ($1 \leq t \leq T$), the processor takes the high-dimensional embedding vectors for node and edge features as inputs and then outputs the representations for the next execution step.

Each node feature vector x_i is projected into *query* (Q_i), *key* (K_i), and *value* (V_i) vectors via learnable parameter matrices W_Q , W_K , and W_V , respectively. Edge features e_{ij} are projected into *key* (K_{ij}) vector with a matrix W_K^E . Then, for each directed edge from node j to node i in the graph G , we compute the attention coefficient

$$\alpha_{ij} = \frac{\langle Q_j, K_i + K_{ij} \rangle}{\sqrt{h}},$$

where $\langle a, b \rangle$ denotes the dot product. Then, each node i normalizes all attention coefficients across its neighbors with the softmax function and temperature τ and receives the aggregated message:

$$\hat{\alpha}_{ij} = \frac{\exp(\alpha_{ij}/\tau)}{\sum_{k \in \mathcal{N}(i)} \exp(\alpha_{ik}/\tau)}, \quad M_i = \sum_{k \in \mathcal{N}(i)} \hat{\alpha}_{ik} V_k, \quad (1)$$

where $\mathcal{N}(i)$ denotes the set of all incoming neighbors of node i and M_i is the message sent to the i -th node.

For undirected graphs, we consider two separate edges in each direction. Also, for each node, we consider a self-loop connecting the node to itself. For multi-head attention, each head l separately computes the messages M_i^l which are then concatenated.

Similar to Transformer Programs, we enforce attention to be hard attention. We found this property important not only for interpretability but also for size generalization, as hard attention allows us to overcome the annealing of the attention weights for arbitrarily large graphs and strictly limits the set of messages that each vertex can receive.

After message computation, node and edge features are updated depending on the current values and sent messages using feed-forward MLP blocks:

$$\begin{aligned} \hat{x}_i^{t+1} &= \text{FFN}_{nodes}([x_i^t, M_i^t]), \\ \hat{e}_{ij}^{t+1} &= \text{FFN}_{edges}([e_{ij}^t, \hat{\alpha}_{ji} V_i, \hat{\alpha}_{ij} V_j]). \end{aligned}$$

We also enforce all node and edge features to be from a fixed finite set, which we call *states*. We ensure such property by adding discrete bottlenecks at the end of the processor block:

$$\begin{aligned} x_i^{t+1} &= \text{Discretize}_{nodes}(\hat{x}_i^{t+1}), \\ e_{ij}^{t+1} &= \text{Discretize}_{edges}(\hat{e}_{ij}^{t+1}). \end{aligned}$$

We implement discretization by projecting the features to the vectors of size k which we force to be one-hot using annealing Gumbel-Softmax (Jang et al., 2017) during training and the argmax at the inference.

3.3 CONTINUOUS INPUTS

Clearly, most of the algorithmic problems operate with continuous or unbounded inputs (e.g., weights on edges). Usually, all input data is encoded into node and edge features and the processor operates over the resulting vectors. The proposed discretization of such vectors would lead to

the loss of information necessary for performing correct execution steps. One possible option to operate with such inputs (we will call them *scalars*, meaning both continuous or size-dependent integer inputs, such as node indexes) is Neural Execution Engines (Yan et al., 2020), which allows one to operate with bit-wise representations of integer and (in theory) real numbers. Such representations are bounded by design, but fully discrete and interpretable.

We propose another option: to maintain scalar inputs (denote them by S) separately from the node and edge features and use them only as edge priorities s_{ij} in the attention block. If scalars are related to the nodes, we assign them to edges depending on the scalar of the sender or receiver node. Now, we can consider the hard attention block as a selector which for each node selects the best edge based on an ordered pair of ‘states priority’ (attention weights described above, which depend only on states of the corresponding nodes and edges) and s_{ij} . We note that this selector is related to the theoretical primitive *select_best* from RASP. We implement this simply by augmenting key vectors K_{ij} of each edge with the indicator if the given edge has the “best” (min or max) scalar among the other edges to node j . Thus, scalars affect only the attention weights, not the messages and the node states.

For multiple different scalar inputs (e.g., weighted edges and node indexes), we use multi-head attention, where each head operates with separate scalars.

Given that, the interface of the proposed processor can be described as

$$X^{t+1}, E^{t+1} = \text{Processor}(X^t, E^t, A, S),$$

where X^t, X^{t+1} and E^t, E^{t+1} are from the fixed sets. State sets are independent of the execution step t and the input graph (including scalar inputs S).

3.4 MANIPULATIONS OVER CONTINUOUS INPUTS

The proposed selector offers a read-only interface to scalar inputs, which is not expressive enough for most of the algorithms. However, we note that the algorithms can be described as *discrete* manipulations over input data. For example, the Dijkstra algorithm (Dijkstra, 1959) takes edge weights as inputs and uses them to find the shortest path distances. Computed distances can affect the consequent execution steps. We note that such distances can be described as the sum of the weights of the edges that form the shortest path to the given node. In other words, the produced scalars depend only on input scalars and discrete execution states.

To avoid the challenges of learning continuous updates with high precision (Klindt, 2023), we propose to learn discrete manipulations with scalars. The updated scalars can then be used with the described selector in the next steps.

In our experiments, we use a scalar updater capable of incrementing, moving, and adding scalars depending on discrete node/edge states:

$$s_i^{t+1} = \text{inc}(x_i^t) + \text{keep}(x_i^t) \cdot s_i^t + \sum_{j \in \mathcal{N}(i)} \text{push}(e_{ji}^t) \cdot s_{ji}^t,$$

$$s_{ij}^{t+1} = \text{inc}(e_{ij}^t) + \text{keep}(e_{ij}^t) \cdot s_{ij}^t + \text{push}(x_i^t) \cdot s_i^t,$$

where s_i are node-related scalars, s_{ij} are edge-related scalars, and *inc*, *keep*, *push* are 0-1 functions representing if scalar in each node/edge should be incremented, kept, or pushed to any of its neighbors. We implement these functions as simple linear projections of node/edge features with consecutive discretization.

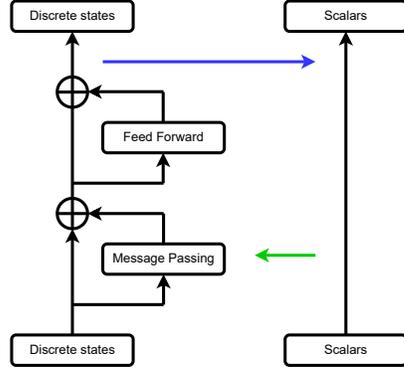


Figure 1: An illustration of the proposed separation between discrete and continuous data flows. Scalars can only affect the attention weights (Green) and can be modified with actions via ScalarUpdate (Blue).

270 Finally, our proposed method (see Figure 1) can be described as:

$$271 \quad X^{t+1}, E^{t+1} = \text{Processor}(X^t, E^t, A, S^t),$$

$$272 \quad S^{t+1} = \text{ScalarUpdate}(X^{t+1}, E^{t+1}, A, S^t).$$

273
274 The proposed neural reasoners are fully discrete and can be interpreted by design. Moreover, the
275 proposed selector block guarantees the predicted behavior of the message passing for any graph
276 size, as it compares discrete state importance and uses continuous scalars only to break ties between
277 equally important states.

278 4 EXPERIMENTS

279 In this section, we perform experiments to evaluate how the proposed discretization affects the
280 performance of neural reasoners on diverse algorithmic reasoning tasks. Our main questions are:

- 281 1. Can the proposed discrete neural reasoners capture the desired algorithmic dynamics with
- 282 hint supervision?
- 283 2. How does discretization affect OOD and size-generalization performance of neural reason-
- 284 ers?
- 285 3. Is the proposed model capable of multi-task learning?

286
287
288 Also, we are interested if discrete neural reasoners can be learned without hints and how they tend
289 to utilize the given amount of node and edge states to solve the problem. We discuss no-hint experi-
290 ments separately in Section 6.

291 4.1 DATASETS

292
293 We perform our experiments on the problems from the recently proposed SALSA-CLRS bench-
294 mark (Minder et al., 2023), namely BFS, DFS, Prim, Dijkstra, Maximum Independent Set (MIS),
295 and Eccentricity. We believe that the proposed method is not limited by the covered problems, but
296 we leave the implementation of the required data flows (e.g., edge-based reasoning (Ibarz et al.,
297 2022), graph-level hints, interactions between different scalars) for future work.

298
299 The train dataset of SALSA-CLRS consists of random graphs with at most 16 nodes sampled from
300 the Erdős-Rényi (ER) distribution with parameter p chosen to be as low as possible while graphs
301 remain connected with high probability. The test set consists of sparse graphs of sizes from 16 to
302 1600 nodes.

303
304 We slightly modify the hints from the benchmark without conceptual changes (e.g., we have modi-
305 fied the hints for the DFS problem to remove graph-level hints). Discrete states are fully described
306 by the non-scalar hints and scalars are exactly the hints of the *scalar* type (we refer to Veličković
307 et al. (2022) for the details on hint design).

308 4.2 BASELINES AND EVALUATION

309
310 We compare the performance of our proposed discrete model with two baseline sparse models, GIN
311 (Xu et al., 2019) and Pointer Graph Network (Veličković et al., 2020a). We report both node-level
312 and graph-level metrics for the baselines and our model. Also, we compare our model with Triplet-
313 GMPNN (Ibarz et al., 2022) and two recent approaches, namely Hint-ReLIC (Bevilacqua et al.,
314 2023) and G-ForgetNet (Bohde et al., 2024), which demonstrate state-of-the-art performance of
315 hint-based neural algorithmic reasoning. However, as these methods are evaluated on the CLRS-30
316 benchmark and their code is not yet publicly available, we can only compare them on the corre-
317 sponding tasks (BFS, DFS, Dijkstra, Prim) and CLRS-30 test data, namely ER graphs with $p = 0.5$
318 of size 64. Note that this test data is more dense than that of SALSA-CLRS, meaning shorter roll-
319 outs for given tasks. Also, only node-level metrics have been reported for these methods.

320 4.3 MODEL DETAILS

321
322 For our experiments, we use the model described in Section 3. We use one attention head for each
323 scalar value. The number of processor steps is defined externally as the length of the ground truth

algorithm trajectory, which is consistent with the prior work. We use one architecture (except the task-dependent encoders/decoders), including the ScalarUpdate module for all the problems.

We recall that neural reasoners can operate either on the base graph (which is defined by the problem) or on more dense graphs with the original graph encoded into edge features. SALSA-CLRS proposes to enhance the size-generalization abilities of neural reasoners with increased test sizes (up to 100-fold train-to-test scaling), so we use the base graph for message-passing, similar to the SALSA-CLRS baselines. We also add a virtual node that communicates with all nodes of the graph.

4.4 TRAINING DETAILS AND HYPERPARAMETERS

We train each model using the Adam optimizer, learning rate $\eta = 0.001$, using teacher forcing, batch size of 32 graphs with 1000 optimization steps, and evaluate the resulting model. We anneal softmax temperatures in the discrete bottlenecks (attention weights, ScalarUpdate operations) geometrically from 3.0 to 0.01, decreasing the temperature at each training step. We report all hyperparameters in the source code.¹

During training, we minimize the standard hints and output losses: scalar hints are optimized with MSE loss, and other types of hints are optimized with cross-entropy and categorized cross-entropy losses (Veličković et al., 2022; Ibarz et al., 2022). Note that we do not supervise any additional details in model behavior, e.g., selecting the most important neighbor in the attention block, the exact operations with scalars, etc.

For multitask experiments, we follow the setup proposed by Ibarz et al. (2022) and train a single processor with task-dependent encoders/decoders to imitate all covered algorithms simultaneously. We make 10000 optimization steps on the accumulated gradients across each task and keep all hyperparameters the same as in the single task.

Our models are trained on a single A100 GPU, requiring less than 1 hour for single-task and 5-6 hours for multitask training.

4.5 RESULTS

We found learning with teacher forcing suitable for discrete neural reasoners, as discretization blocks allow us to perform the exact transitions between the states. Trained with step-wise hint supervision, discrete neural reasoners are able to perfectly align with the original algorithm and generalize on larger test data without any performance loss. We report the evaluation results in Tables 1 and 2. Also, our multitask experiments show that the proposed discrete models are capable of multitask learning and demonstrate the perfect generalization scores in a multitask manner too.

Recall that we have three key components of our contribution: feature discretization, hard attention, and separating discrete and continuous data flows. To evaluate the importance of each component for generalization capabilities of the proposed models, we conduct an ablation study, the details can be found in Appendix A. In

Table 1: Node \ graph level test scores for the proposed discrete reasoner and the baselines on SALSA-CLRS test data. Scores are averaged across 5 different seeds, standard deviation is omitted. For the eccentricity problem, only the graph-level metric is applicable.

TASK	SIZE	GIN	PGN	DNAR (OURS)
BFS	16	98.8 \ 92.5	100. \ 100.	100. \ 100.
	80	95.3 \ 59.4	99.8 \ 88.1	100. \ 100.
	160	95.1 \ 37.8	99.6 \ 66.3	100. \ 100.
	800	86.9 \ 0.9	98.7 \ 0.2	100. \ 100.
	1600	86.5 \ 0.0	98.5 \ 0.0	100. \ 100.
DFS	16	41.5 \ 0.0	82.0 \ 19.9	100. \ 100.
	80	30.4 \ 0.0	38.4 \ 0.0	100. \ 100.
	160	20.0 \ 0.0	26.9 \ 0.0	100. \ 100.
	800	19.5 \ 0.0	24.9 \ 0.0	100. \ 100.
	1600	17.8 \ 0.0	23.1 \ 0.0	100. \ 100.
SP	16	95.2 \ 49.8	99.3 \ 89.5	100. \ 100.
	80	62.4 \ 0.0	94.2 \ 3.3	100. \ 100.
	160	53.3 \ 0.0	92.0 \ 0.0	100. \ 100.
	800	40.4 \ 0.0	87.1 \ 0.0	100. \ 100.
	1600	36.9 \ 0.0	84.5 \ 0.0	100. \ 100.
PRIM	16	89.6 \ 29.7	96.4 \ 69.9	100. \ 100.
	80	51.6 \ 0.0	79.7 \ 0.0	100. \ 100.
	160	49.5 \ 0.0	75.6 \ 0.0	100. \ 100.
	800	45.0 \ 0.0	69.5 \ 0.0	100. \ 100.
	1600	43.2 \ 0.0	66.8 \ 0.0	100. \ 100.
MIS	16	79.9 \ 3.3	99.8 \ 98.6	100. \ 100.
	80	79.9 \ 20.0	99.4 \ 88.9	100. \ 100.
	160	78.2 \ 0.0	99.4 \ 76.2	100. \ 100.
	800	83.4 \ 0.0	98.8 \ 18.0	100. \ 100.
	1600	79.2 \ 0.0	98.9 \ 5.2	100. \ 100.
Ecc.	16	25.3	100.	100.
	80	23.8	100.	100.
	160	26.1	100.	100.
	800	17.1	100.	100.
	1600	16.0	83.0	100.

¹<https://anonymous.4open.science/r/F4CA/>

Table 2: Node-level test scores for the proposed discrete reasoner and the baselines on CLRS-30 test data. Test graphs are of size 64. Scores are averaged across 5 different seeds.

TASK	TRIPLET-GMPNN	HINT-RELIC	G-FORGETNET	DNAR (OURS)
BFS	99.73 \pm 0.0	99.00 \pm 0.2	99.96 \pm 0.0	100. \pm 0.0
DFS	47.79 \pm 4.2	100. \pm 0.0	74.31 \pm 5.0	100. \pm 0.0
DIJKSTRA	96.05 \pm 0.6	97.74 \pm 0.5	99.14 \pm 0.1	100. \pm 0.0
MST-PRIM	86.39 \pm 1.3	87.97 \pm 2.9	95.19 \pm 0.3	100. \pm 0.0

short, our additional experiments demonstrate that the proposed processor without the discretization performs non-above the baselines level and that removing hard-attention or discrete *ScalarUpdate* module strictly limits the generalization capabilities of the proposed model.

5 INTERPETABILITY AND TESTING

In addition to the empirical evaluation of the trained models on diverse test data, the discrete and size-independent design of the proposed models allows us to interpret and test them manually. The main idea is to show that the model will perform the exact discrete state transitions (including discrete operations with scalars) as the original algorithm.

First, we note that due to the hard attention, each node receives exactly one message. Also, the message depends only on the discrete states of the corresponding nodes and edges. Thereby, as each node and edge states after a single processor step depend only on the current states and received message, all the possible options can be directly enumerated and tested if all states change to the correct ones.

The only remaining part to fully interpret the whole model is the attention block. We note that our implementation of the *select_best* selector (Section 3.3) does not necessarily produce the top-1 choice over the ordered pairs of ‘states priority’ and scalars s_{ij} as it simply augments key vectors with indicators if the given edge has the best scalar among others. For example, it may happen that for some state, the maximum attention weight is achieved for an edge without the indicator. However, given the finite number of discrete states, we can manually check if the mentioned “best” indicator increases the attention weight between every pair of states. Combined with hard attention, this would imply that the attention block attends depending on the predefined states and uses scalar priorities only to break ties. Please see the particular example with more detailed explanation for the BFS problem in the Appendix D.

Given that, we can unit-test all possible state transitions and attention blocks. With full coverage of such tests, we can guarantee the correct execution of the desired algorithm for *any* test data. We tested our trained models from Section 4 manually verifying state transitions. As a result, we confirm that the attention block indeed operates as *select_best* selector, as the model actually uses these indicators to increase the attention weights. Thus, we can guarantee that for any graph size, the model will mirror the desired algorithm, which is correct for any test size.

6 TOWARDS NO-HINT DISCRETE REASONERS

In this section, we discuss the challenges of training discrete reasoners without hints, which can be useful when tackling new algorithmic problems.

Training deep discretized models is known to be challenging: without hyperparameter search, discrete models are only slightly improved over the untrained models. Therefore, we focus only on the BFS algorithm, as it is well-aligned with the message-passing framework, has short roll-outs, and can be solved with small states count (note that for no-hint models node/edge states count is a hyperparameter due to the absence of the ground truth states trajectory).

We recall that the output of the BFS problem is the exploration tree pointing from each node to its parent. Each node chooses as a parent the neighbor from the previous distance layer with the smallest index.

Table 3: BFS node/graph level scores of the *best_no_hint_model* for different graph sizes.

	5	16	64
<i>best_no_hint_model</i>	99 / 86	94 / 34	79 / 0

We perform hyperparameter search over the training sizes (using ER graphs with $p = 0.5$ and $n \in [4, 16]$), discrete node states count (from 2 to 6 states), softmax temperature annealing schedules ($[3, 0.01]$, $[3, 0.1]$, $[3, 1]$). For each hyperparameter choice, we train 5 models with different seeds. We validate the resulting models on the graphs of size 16. The best resulting model is obtained with the training size 5 and 4 node states. The trained models never achieved the perfect validation scores, see Table 3 for the results.

Then, we select the best-performing models and try to analyze the mistakes of the resulting models and reverse-engineer how they utilize the given states. First, we look at the node states after the last step of the processor and note that the states correspond to the distances from the starting node. More formally, we note that the model with four states uses the first state for the starting node, the second state for its neighbors, the third state for nodes at distance two from the starting node, and the last state for all other nodes and such states-based classification of distance has accuracy $> 98\%$ when tested on 1000 random graphs with 16 nodes. Then, we note that for the nodes that are from the first four distance layers from the starting node, the pointers are predicted with 100% accuracy and these pointers are computed layer-by-layer as in the ground truth algorithm (we refer to Appendix B for illustrations). The mistakes of the model are on the distance ≥ 4 from the starting node (we did not reverse-engineer the specific logic of computations on larger distances).

We found this behavior well-aligned to the BFS algorithm and indicating the possibility of achieving the perfect validation score with enough states count. However, this algorithm does not generalize since it fails at distances larger than those encountered during training.

On the other hand, one can demonstrate that for a small enough state count (for BFS, it is two node states and two edge states) and diverse enough validation data, the perfect validation performance implies that the learned solution will generalize to any graph size.

Therefore, we highlight the need to achieve perfect validation performance with models that use as few states as possible, which corresponds to the minimum description length (MDL) theory (Myung, 2000; Rissanen, 2006) and is related to the notion of Kolmogorov complexity (Kolmogorov, 1963).

Finally, we note that for sequential problems, such as DFS, obtaining a good no-hint model can be even more challenging and can require additional effort. One possible way to overcome this limitation is to implement a curriculum learning setup.

7 LIMITATIONS AND FUTURE WORK

Limitations In this work, we propose a method to learn robust neural reasoners that demonstrate perfect generalization performance and are interpretable by design. In this section, we describe some limitations of our work and important directions for future research.

First, several proposed design choices strictly reduce the expressive power of the model. For example, due to the hard attention, the proposed model is unable to compute the mean value from all the neighbors in a single message-passing step, which is trivial for attention-based models (note that this can be computed in several message-passing steps). Thus, the model in its current form is unable to express transitions between hints for some algorithms from the CLRS-30 in a single processor step. However, we believe that the expressivity of the proposed model can be enhanced with additional architectural modifications (e.g., edge-based reasoning (Ibarz et al., 2022), global states, interactions between different scalars) that can be combined with the proposed discretization ideas.

Second, while we report the perfect scores for the covered tasks, we cannot guarantee that the training will converge to the correct model for any initialization/training data distributions. However, we empirically found the proposed method to be quite robust to various architecture/training hyperparameter choices.

Future work Our method is based on the particular architectural choice and actively utilizes the attention mechanism. However, the graph deep learning field is rich in various architectures exploiting different inductive biases and computation flows. The proposed separation between discrete states and continuous inputs may apply to other models, however, any particular construction can require additional efforts.

Also, we provide only one example of the `ScalarUpdate` block. We believe that utilizing a general architecture (e.g., some form of discrete Neural Turing Machine (Graves et al., 2014; Gulcehre et al., 2016)) capable of executing a wider range of manipulation is of interest for future work.

With the development of neural reasoners and their ability to execute classic algorithms on abstract data, it is becoming more important to investigate how such models can be applicable in real-world scenarios according to the Neural Algorithmic Reasoning blueprint (Veličković & Blundell, 2021) and transfer their knowledge to high-dimensional noisy data with intrinsic algorithmic dependencies. While there are several examples of NAR-based models tackling real-world problems (Beurer-Kellner et al., 2022; Numeroso et al., 2023), there are no established benchmarks for extensive evaluation and comparison of different approaches.

Lastly, we leave for future work a deeper investigation of learning interpretable neural reasoners without hints, which we consider essential from both theoretical perspective and practical applications, e.g., combinatorial optimization.

8 CONCLUSION

In this paper, we force neural reasoners to maintain the execution trajectory as a combination of finite predefined states. To achieve that, we separate discrete and continuous data flows and describe the interaction between them. The obtained discrete reasoners are interpretable by design. Moreover, trained with hint supervision, such models perfectly capture the dynamic of the underlying algorithms and do not suffer from distributional shifts. We consider discretization of hidden representations as a crucial component for robust neural reasoners.

REFERENCES

- Luca Beurer-Kellner, Martin Vechev, Laurent Vanbever, and Petar Veličković. Learning to configure computer networks with neural algorithmic reasoning. *Advances in Neural Information Processing Systems*, 35:730–742, 2022.
- Beatrice Bevilacqua, Kyriacos Nikiforou, Borja Ibarz, Ioana Bica, Michela Paganini, Charles Blundell, Jovana Mitrovic, and Petar Veličković. Neural algorithmic reasoning with causal regularisation. In *International Conference on Machine Learning*, 2023.
- Montgomery Bohde, Meng Liu, Alexandra Saxton, and Shuiwang Ji. On the markov property of neural algorithmic reasoning: Analyses and methods. In *The Twelfth International Conference on Learning Representations*, 2024.
- Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms*. MIT press, 2009.
- Artur Back de Luca and Kimon Fountoulakis. Simulation of graph algorithms with looped transformers. 2024.
- Edsger W. Dijkstra. A note on two problems in connexion with graphs. *Numerische Mathematik*, 1: 269–271, 1959.
- Andrew J Dudzik and Petar Veličković. Graph neural networks are dynamic programmers. *Advances in Neural Information Processing Systems*, 35:20635–20647, 2022.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and

- 540 Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*,
541 2021. <https://transformer-circuits.pub/2021/framework/index.html>.
- 542 Dan Friedman, Alexander Wettig, and Danqi Chen. Learning transformer programs. *Advances in*
543 *Neural Information Processing Systems*, 2023.
- 544 Dobrik Georgiev, Pietro Barbiero, Dmitry Kazhdan, Petar Velivckovi'c, and Pietro Lio'. Algorith-
545 mic concept-based explainable reasoning. In *AAAI Conference on Artificial Intelligence*, 2021.
546 URL <https://api.semanticscholar.org/CorpusID:235899244>.
- 547 Dobrik Georgiev, Pietro Lio, Jakub Bachurski, Junhua Chen, and Tunan Shi. Beyond erdos-renyi:
548 Generalization in algorithmic reasoning on graphs. In *NeurIPS 2023 Workshop: New Frontiers*
549 *in Graph Learning*, 2023.
- 550 Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines. *arXiv preprint*
551 *arXiv:1410.5401*, 2014.
- 552 Caglar Gulcehre, Sarath Chandar, Kyunghyun Cho, and Yoshua Bengio. Dynamic neural turing
553 machine with soft and hard addressing schemes. *arXiv preprint arXiv:1607.00036*, 2016.
- 554 Jessica B Hamrick, Kelsey R Allen, Victor Bapst, Tina Zhu, Kevin R McKee, Joshua B Tenen-
555 baum, and Peter W Battaglia. Relational inductive bias for physical construction in humans and
556 machines. *arXiv preprint arXiv:1806.01203*, 2018.
- 557 Borja Ibarz, Vitaly Kurin, George Papamakarios, Kyriacos Nikiforou, Mehdi Bennani, Róbert
558 Csordás, Andrew Joseph Dudzik, Matko Bošnjak, Alex Vitvitskyi, Yulia Rubanova, et al. A
559 generalist neural algorithmic learner. In *Learning on Graphs Conference*. PMLR, 2022.
- 560 Rishabh Jain, Petar Veličković, and Pietro Liò. Neural priority queues for graph neural networks.
561 *arXiv preprint arXiv:2307.09660*, 2023.
- 562 Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *In-*
563 *ternational Conference on Learning Representations*, 2017.
- 564 Jonas Jürß, Dulhan Hansaja Jayalath, and Petar Veličković. Recursive algorithmic reasoning. pp.
565 5–1, 2024.
- 566 Łukasz Kaiser and Ilya Sutskever. Neural GPUs learn algorithms. *arXiv preprint arXiv:1511.08228*,
567 2015.
- 568 David A. Klindt. Controlling neural network smoothness for neural algorithmic reasoning. *Trans-*
569 *actions on Machine Learning Research*, 2023. ISSN 2835-8856.
- 570 Andrei N. Kolmogorov. On tables of random numbers. *Sankhya: The Indian Journal of Statistics*,
571 25, 1963.
- 572 David Lindner, János Kramár, Sebastian Farquhar, Matthew Rahtz, Tom McGrath, and Vladimir
573 Mikulik. Tracr: Compiled transformers as a laboratory for interpretability. *Advances in Neural*
574 *Information Processing Systems*, 36, 2024.
- 575 Sadeh Mahdavi, Kevin Swersky, Thomas Kipf, Milad Hashemi, Christos Thrampoulidis, and Ren-
576 jie Liao. Towards better out-of-distribution generalization of neural algorithmic reasoning tasks.
577 *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856.
- 578 Julian Minder, Florian Grötschla, Joël Mathys, and Roger Wattenhofer. SALSA-CLRS: A sparse
579 and scalable benchmark for algorithmic reasoning. *arXiv preprint arXiv:2309.12253*, 2023.
- 580 In Jae Myung. The importance of complexity in model selection. *Journal of mathematical psychol-*
581 *ogy*, 44(1):190–204, 2000.
- 582 Danilo Numeroso, Davide Bacciu, and Petar Veličković. Dual algorithmic reasoning. In *Internat-*
583 *ional Conference on Learning Representations*, 2023.
- 584 Scott Reed and Nando De Freitas. Neural programmer-interpreters. *arXiv preprint*
585 *arXiv:1511.06279*, 2015.

- 594 Jorma Rissanen. Information and complexity in statistical modeling. In *Information Theory Work-*
595 *shop*, 2006.
- 596
- 597 Gleb Rodionov and Liudmila Prokhorenkova. Neural algorithmic reasoning without intermediate
598 supervision. *Advances in Neural Information Processing Systems*, 2023.
- 599
- 600 Dan Roth Roni Khardon. Learning to reason. In *Proceedings of the 12th National Conference on*
601 *Artificial Intelligence, Seattle, WA, USA, July 31 - August 4, 1994, Volume 1*, pp. 682–687. AAAI
602 Press / The MIT Press, 1994.
- 603
- 604 Yunsheng Shi, Zhengjie Huang, Shikun Feng, Hui Zhong, Wenjin Wang, and Yu Sun. Masked label
605 prediction: Unified message passing model for semi-supervised classification. *arXiv preprint*
arXiv:2009.03509, 2020.
- 606
- 607 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
608 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*
tion processing systems, 30, 2017.
- 609
- 610 Petar Veličković and Charles Blundell. Neural algorithmic reasoning. *Patterns*, 2(7):100273, 2021.
- 611
- 612 Petar Veličković, Lars Buesing, Matthew Overlan, Razvan Pascanu, Oriol Vinyals, and Charles
613 Blundell. Pointer graph networks. *Advances in Neural Information Processing Systems*, 33:
2232–2244, 2020a.
- 614
- 615 Petar Veličković, Rex Ying, Matilde Padovano, Raia Hadsell, and Charles Blundell. Neural execu-
616 tion of graph algorithms. In *International Conference on Learning Representations*, 2020b.
- 617
- 618 Petar Veličković, Adrià Puigdomènech Badia, David Budden, Razvan Pascanu, Andrea Banino,
619 Misha Dashevskiy, Raia Hadsell, and Charles Blundell. The CLRS algorithmic reasoning bench-
620 mark. In *International Conference on Machine Learning*, pp. 22084–22102. PMLR, 2022.
- 621
- 622 Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. *Advances in neural informa-*
tion processing systems, 28, 2015.
- 623
- 624 Gail Weiss, Yoav Goldberg, and Eran Yahav. Thinking like transformers. In *International Confer-*
ence on Machine Learning, pp. 11080–11090. PMLR, 2021.
- 625
- 626 Louis-Pascal Xhonneux, Andreea-Ioana Deac, Petar Veličković, and Jian Tang. How to transfer
627 algorithmic reasoning knowledge to learn new algorithms? *Advances in Neural Information*
Processing Systems, 34:19500–19512, 2021.
- 628
- 629 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
630 networks? In *The Seventh International Conference on Learning Representations*, 2019.
- 631
- 632 Keyulu Xu, Jingling Li, Mozhi Zhang, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie Jegelka.
633 What can neural networks reason about? In *International Conference on Learning Representa-*
tions, 2020.
- 634
- 635 Yujun Yan, Kevin Swersky, Danai Koutra, Parthasarathy Ranganathan, and Milad Hashemi. Neural
636 execution engines: Learning to execute subroutines. *Advances in Neural Information Processing*
Systems, 33:17298–17308, 2020.
- 637
- 638 Wojciech Zaremba and Ilya Sutskever. Learning to execute. *arXiv preprint arXiv:1410.4615*, 2014.
- 639
- 640 Hattie Zhou, Arwen Bradley, Etai Littwin, Noam Razin, Omid Saremi, Joshua M Susskind, Samy
641 Bengio, and Preetum Nakkiran. What algorithms can transformers learn? A study in length
642 generalization. In *The Twelfth International Conference on Learning Representations*, 2024.
- 643
- 644
- 645
- 646
- 647

A ABLATION STUDY

Recall that we have three key components of our contribution: feature discretization, hard attention, and separating discrete and continuous data flows. In this section, we study the importance of these components for generalization capabilities of the proposed models.

Discrete bottlenecks First, we evaluate the model without all discrete bottlenecks: the result is a simple Transformer Convolution processor, which performs comparable to other baseline models, see Table 4.

Table 4: Node \ graph level test scores for our base model without all discrete bottlenecks. Scores are averaged across 5 different seeds, standard deviation is omitted.

SIZE	16	80	160	800	1600
BFS	99.9 \ 99.3	99.7 \ 88.2	99.5 \ 57.9	98.4 \ 0.0	97.2 \ 0.0
DFS	79.2 \ 6.8	41.1 \ 0.0	28.1 \ 0.0	24.7 \ 0.0	21.9 \ 0.0
SP	99.3 \ 88.7	94.1 \ 12.4	90.3 \ 0.0	86.9 \ 0.0	82.4 \ 0.0
PRIM	95.1 \ 72.7	82.6 \ 0.0	79.7 \ 0.0	68.1 \ 0.0	66.0 \ 0.0
MIS	99.8 \ 98.6	99.6 \ 86.1	99.2 \ 69.0	97.1 \ 11.9	96.3 \ 0.0
ECC	79.2	41.1	28.1	24.7	21.9

Hard attention To highlight the importance of the hard attention for strong size generalization, we train the proposed model but with the regular attention mechanism on the BFS task. The resulting model also demonstrates the perfect scores for the given test data. However, the standard test data (the Erdős-Renyi graphs with low edge probability) does not contain nodes with a large number of neighbors, while such nodes can be problematic due to the annealing of the attention weights. Thus, we additionally test the resulting model on the complete bipartite graphs $K_{2,n-2}$ for different n . For each n , we assign the starting node from the smaller component and test if the second node in this component correctly selects its parent in the BFS tree, see Figure 2.

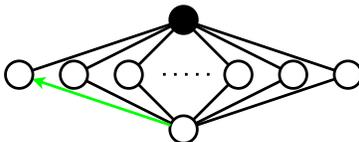


Figure 2: Complete bipartite graphs $K_{2,n-2}$ used to evaluate the effect of the attention weights annealing. The black node is the starting node. The highlighted edge (Green) is the ground truth pointer from the bottom node to its parent from the BFS tree.

Our experiments demonstrate that the model without hard attention fails to predict the correct pointer for larger graphs due to the attention weight annealing, see Table 5. We note that the models from Section 4 are provably correct on any test data.

Table 5: Attention weights for the ground truth pointer (green pointer from Figure 2) for different graph sizes; (+/-) denotes if the correct pointer was predicted.

SIZE	16	80	160	800	1600
ATTENTION WEIGHT	0.97 (+)	0.86 (+)	0.76 (+)	0.38 (-)	0.24 (-)

Scalar updater As demonstrated in Klindt (2023), simple neural networks trained to sum two real numbers fail to learn the structure of the task and struggle to extrapolate beyond the training data distributions. In this section, we study how these limitations affect the overall performance of neural reasoners to highlight the importance of the proposed discrete manipulations with scalars.

First, we recall that usually all input data is encoded into node and edge features and the processor operates over the resulting vectors. Then, hints of type scalar are directly predicted from the

node/edge features. To evaluate the effect of non-discrete *ScalarUpdate* modules, we simply replace the proposed discrete *ScalarUpdate* module with a single-layer transformer convolution network, which inputs scalars and node/edge states and outputs the scalars for the next step, keeping the remaining architecture the same as in the main experiments. We train the resulting model on Dijkstra and MST problems.

Additionally, we evaluate the non-discrete *ScalarUpdate* module in a more straightforward setup. Similarly to Klindt (2023), we train a 2-layer MLP to add two real numbers and use the resulting model as a *ScalarUpdate* module for the Dijkstra algorithm. We manually use the learned addition module when the node distances are updated (e.g., the distance of the node u is updated with the sum of the distance of v and the edge (v, u) cost), and use the ground truth scalars for other scalar updates. Our experiments demonstrate that the resulting model outperforms the baselines on the test size of 16 nodes, but does not generalize well on larger graphs, see the evaluation results in Table 6.

Table 6: Node \ graph level test scores for the proposed model with *ScalarUpdater* replaced by a regular attention-based network trained to predict hints of type scalar. ‘Addition only’ means that *ScalarUpdater* is replaced by a 2-layer MLP trained to predict the sum of two numbers (other values are taken from the ground truth).

SIZE	16	80	160	800	1600
DIJKSTRA	99.3 \ 94.6	60.7 \ 0.0	42.8 \ 0.0	19.0 \ 0.0	11.8 \ 0.0
MST	99.8 \ 98.1	98.2 \ 54.1	97.2 \ 28.1	95.5 \ 0.0	91.73 \ 0.0
DIJKSTRA (ADDITION ONLY)	99.8 \ 96.6	95.3 \ 71.0	86.5 \ 46.3	41.6 \ 3.1	22.2 \ 0.0

We note that all the resulting models demonstrate perfect scores when evaluated with teacher-forced ground truth scalars. Thus, all state transitions are learned correctly and imperfect test scores are fully described by the errors in manipulations with continuous values.

To summarize, small errors in manipulations with scalars (even restricted on the simplest addition sub-task) strictly affect the overall performance of the model, highlighting the importance of the proposed discrete manipulations with scalars.

B STATE USAGE FOR NO-HINT MODELS

In this section, we provide several illustrations of the node states and the dynamics of the pointer prediction updates (Figures 3-5). Our analysis suggests that no-hint models with K states tend to use states as distances from the starting node, with the distances $\geq K$ merged to the same state. Also, pointer predictions for the first K BFS layers are correct and computed layer-by-layer as in the ground truth algorithm. The mistakes of the model are at the later layers and some pointers at the later layers are computed before that in the ground truth algorithm (Figure 5c). For the simplicity of illustrations, we use the model with 3 discrete states.

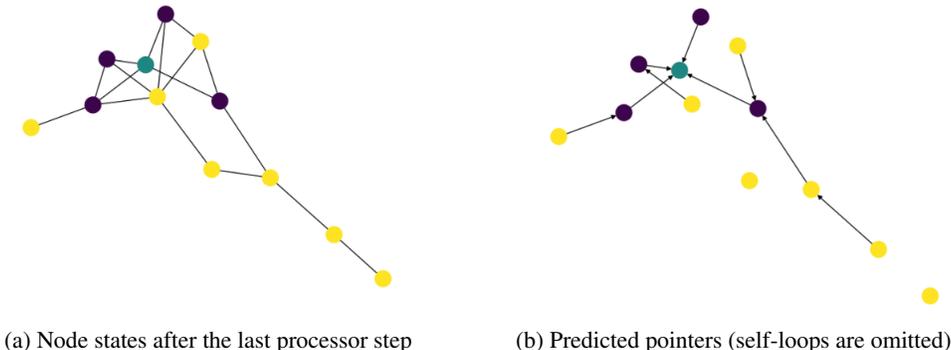
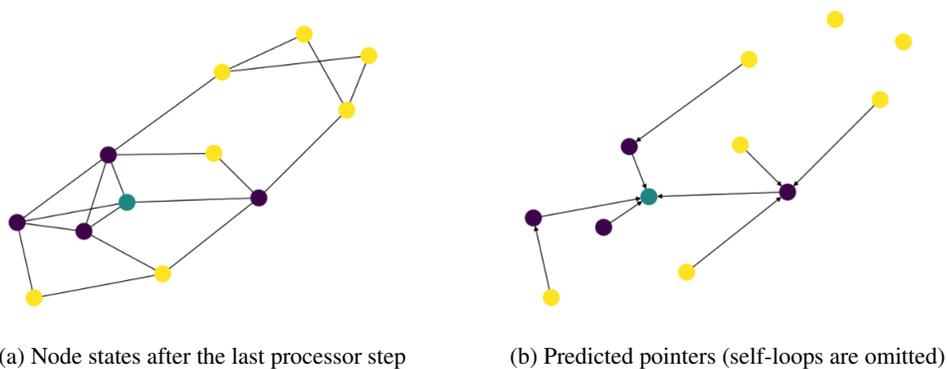


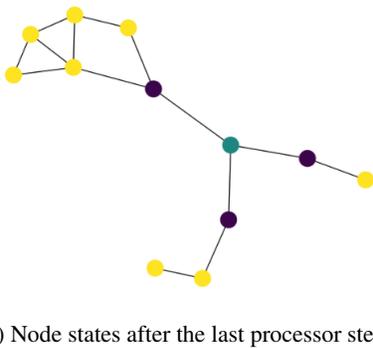
Figure 3: Node states and the predicted pointers after the last processor step of the DNAR model (with 3 states), trained without hints. Different colors represent different states. The green node is the starting node.

756
757
758
759
760
761
762
763
764
765
766
767

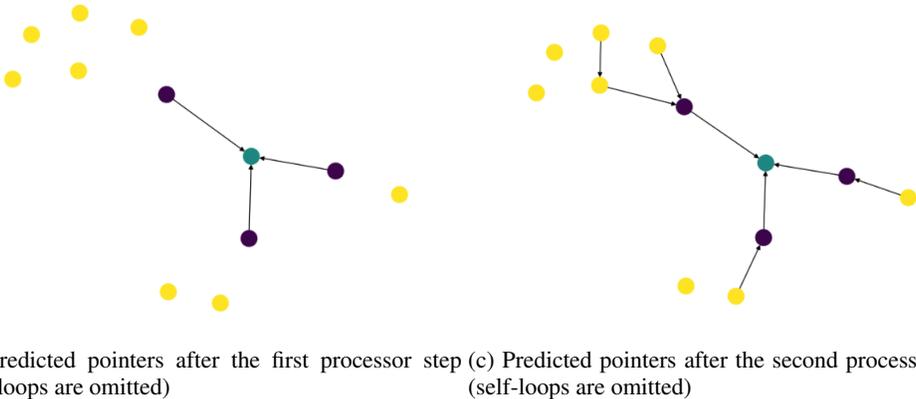


768 Figure 4: Node states and the predicted pointers after the last processor step of the DNAR model
769 (with 3 states), trained without hints. Different colors represent different states. The green node is
770 the starting node.
771

772
773
774
775
776
777
778
779
780
781
782
783



784
785
786
787
788
789
790
791
792
793
794
795
796
797



798 Figure 5: Node states and the dynamics of the pointer prediction updates of the DNAR model (with
799 3 states), trained without hints. Different colors represent different states. The green node is the
800 starting node.
801

802
803 **C EXTENDED *ScalarUpdate* MODULE**

804
805 In this section, we investigate if the proposed *ScalarUpdate* module can be successfully extended
806 to support more complex manipulations with scalars.
807

808 First, we note that simple manipulations with scalars cover a significant part of the classical al-
809 gorithms. In this work, we use the minimum set of the required functions, but it can be directly
extended by other functions. Importantly, as *ScalarUpdate* can be viewed as a separate module,

we can separately check if it is possible to train it with any given set of predefined manipulations for any problem only with supervision on the results.

Let us formalize the problem: each input of the *ScalarUpdate* module can be described as an object with a discrete state s_i (from a fixed predefined set) and several scalar values (we consider two scalars x_i and y_i). Note that we omit the separation between nodes and edges and consider objects with several scalar values. For each discrete state s_i , there exists a ground truth update of scalars, e.g., $f(s, x, y) = x + \cos(y)$. The output of the scalar updater can be viewed as a sum:

$$\text{ScalarUpdate}(s, x, y) = \sum_{g \in OPS} g(x, y) \cdot \text{activate}_g(s),$$

where OPS is a predefined set of operations and activate_g is a 0-1 function representing if a specific operation should be applied. Note that activate_g depends only on a discrete state of the input.

For our additional experiment, we train several different *ScalarUpdate* modules with the extended set of operations:

- $g_0(x, y) = 1$;
- $g_1(x, y) = x$;
- $g_2(x, y) = \cos(x)$;
- $g_3(x, y) = x \cdot y$;
- $g_4(x, y) = \text{atan2}(x, y)$

to learn the following set of the ground truth updates simultaneously:

- $f_0(x, y) = x$;
- $f_1(x, y) = \cos(x)$;
- $f_2(x, y) = \cos(x) + x \cdot y$;
- $f_3(x, y) = \text{atan2}(x, y)$;
- $f_4(x, y) = 1 + x + \text{atan2}(x, y)$.

In particular, we consider a set of 16 discrete states (numbered from 0 to 15 and sampled uniformly) and the ground truth scalar update is derived from these states by taking the remainder of the division by 5 (updates count).

The learnable parameters of the *ScalarUpdate* are state’s embeddings and linear projections for each indicator. We train *ScalarUpdate* to minimize the MSE loss between the ground truth and predicted outputs with 5000 optimization steps. Additionally, we train a non-discrete scalar updater (2-layer MLP), similar to our ablation experiments. We refer to the source code for the experiment details.

Inspired by Klindt (2023), we generate training scalars X and Y from $\text{Uniform}[0.5, 1.0]$ and generate test set by sampling scalars from the $\text{Uniform}[0., 0.5]$ distribution.

We report the evaluation results in Table 7. The proposed discrete *ScalarUpdate* module successfully learned the correct operations for updates f_1, \dots, f_4 for all seeds and 3 times out of 5 for f_0 (note that the model was trained to predict different manipulations for different states simultaneously). For unsuccessful runs, when f_0 was not learned correctly, the learned operation for f_0 was g_2 for some states (i.e. $x \cdot y$ instead of x), which can be explained by optimization challenges as the distribution of y is close to 1.

Our experiment demonstrates that the proposed *ScalarUpdate* module can be extended to support a wider range of manipulations with scalars. We note that this complicates the optimization problem of selecting the correct operations/operands from the operations results (e.g., such decomposition might not be unique).

D INTERPETABILITY AND TESTING DETAILS

In this section, we provide additional details on interpretability and testing of the proposed discrete reasoners.

Table 7: MSE for train/test distributions for the discrete and non-discrete *ScalarUpdate* modules and different operations.

	f_0	f_1	f_2	f_3	f_4
<i>discrete</i>	0.01 / 0.1	0. / 0.	0 / 0	0 / 0	0 / 0
<i>non-discrete</i>	$5.7e-6 / 0.03$	$5.1e-6 / 0.001$	$1.1e-5 / 0.007$	$1.0e-5 / 0.08$	$2.0e-5 / 0.03$

As an example, consider the BFS algorithm. First, recall the pseudocode of the algorithm:

```

Starting_node ← visited
All_other_nodes ← not_visited
for step in range(T) do
  for node  $U$  in a graph do
    if  $U$  is visited on previous steps then
      continue
    end if
    if  $U$  has a neighbor  $P$  that visited on previous steps then
       $U \leftarrow$  visited on this step
       $U$  select the smallest-indexed such neighbor  $P$  as parent:
      Edge ( $U, P$ ) ← pointer
      Self-loop ( $U, U$ ) ← not_pointer
    end if
  end for
end for
return a BFS tree described by pointers

```

Now let us describe how we can verify that the trained DNAR model will perfectly imitate this algorithm for any test data.

First, we note that for each node U , the node state on the step $t+1$ (denoted by U_{t+1}) is the function of U_t and V_t , where V_t is the node that sends the message to U on step t :

$$U_{t+1} = \text{StateUpdate}(U_t, \text{message_from_}V_t)$$

How does the node U select a node that will send a message to it? For any node V connected to U , the node U computes attention scores depending on discrete states of V of each node and a discrete indicator if each node has the smallest (or largest) scalar among all neighbors of U with the same discrete state as V . Then, the node U selects the node V with the largest attention score.

In our (slightly simplified) case, the attention scores only depend on the tuples ($U_{state}, V_{state}, \text{indicator_if_}u\text{ has the smallest index}$) and there are only 8 such tuples. We can directly compute these attention scores and verify the required invariants, e.g.,

$$\text{Attention}(\text{not_visited}, \text{visited}, \text{smallest}) > \text{Attention}(\text{not_visited}, *any_other*),$$

which would imply that the not_visited node will receive the message from the smallest-indexed visited neighbor if such exists independently of the graph size and distribution. If there is no such neighbor, the node U will receive the message from another not_visited node (or from itself).

After verifying the correctness of the message flows, we need to ensure that the state updates are computed correctly, e.g.,

$$\text{visited} = \text{StateUpdate}(\text{not_visited}, \text{message_from_visited})$$

$$\text{visited} = \text{StateUpdate}(\text{visited}, *any*)$$

$$\text{not_visited} = \text{StateUpdate}(\text{not_visited}, \text{message_from_not_visited})$$

The main idea is that due to the finite states count and discrete manipulations with scalars, there are only finite amounts of such checks that can cover all possible state transitions and all of them should be evaluated only once.