
Algorithmic Concept-based Explainable Reasoning

Dobrik Georgiev
University of Cambridge
dgg30@cam.ac.uk

Pietro Barbiero
University of Cambridge
pb737@cam.ac.uk

Dmitry Kazhdan
University of Cambridge
dk525@cam.ac.uk

Petar Veličković
DeepMind
petarv@deepmind.com

Pietro Liò
University of Cambridge
pl219@cam.ac.uk

Abstract

Recent research on graph neural network (GNN) models successfully applied GNNs to classical graph algorithms and combinatorial optimisation problems. This has numerous benefits, such as allowing applications of algorithms when preconditions are not satisfied, or reusing learned models when sufficient training data is not available or can't be generated. Unfortunately, a key hindrance of these approaches is their lack of explainability, since GNNs are black-box models that cannot be interpreted directly. In this work, we address this limitation by applying existing work on concept-based explanations to GNN models. We introduce *concept-bottleneck GNNs*, which rely on a modification to the GNN readout mechanism. Using three case studies we demonstrate that: (i) our proposed model is capable of accurately learning concepts and extracting propositional formulas based on the learned concepts for each target class; (ii) our concept-based GNN models achieve comparative performance with state-of-the-art models; (iii) we can derive global graph concepts, without explicitly providing any supervision on graph-level concepts.

1 Introduction

Graph neural networks (GNNs) have successfully been applied to problems involving data with irregular structure, such as quantum chemistry (Gilmer et al., 2017), drug discovery (Stokes et al., 2020), social networks (Pal et al., 2020) and physics simulations (Battaglia et al., 2016). One of the latest areas of GNN research focuses on using GNNs for emulation of classical algorithms (Cappart et al., 2021). In particular, this research explored applications of GNNs to iterative algorithms (Veličković et al., 2020b; Georgiev and Liò, 2020), pointer-based data structures (Veličković et al., 2020a; Strathmann et al., 2021), and even planning tasks (Deac et al., 2020). Importantly, these works demonstrate that GNNs are capable of *strongly generalising* to input graphs much larger than the ones seen during training.

Unfortunately, in all of the aforementioned cases, these state-of-the-art GNN models are black-boxes, whose behaviour cannot be understood/interpreted directly. In practice, this can lead to a lack of trust in such models, making it challenging to apply and regulate these models in safety-critical applications, such as healthcare. Furthermore, this lack of interpretability also makes it difficult to extract the knowledge learned by such models, which prevents users from better understanding the corresponding tasks (Adadi and Berrada, 2018; Molnar, 2020; Doshi-Velez and Kim, 2017).

Recent work on Explainable AI (XAI) introduced a novel type of Convolutional Neural Network (CNN) explanation approach, referred to as *concept-based explainability* (Koh et al., 2020; Kazhdan et al., 2020b; Ghorbani et al., 2019; Kazhdan et al., 2021). Concept-based explanation approaches

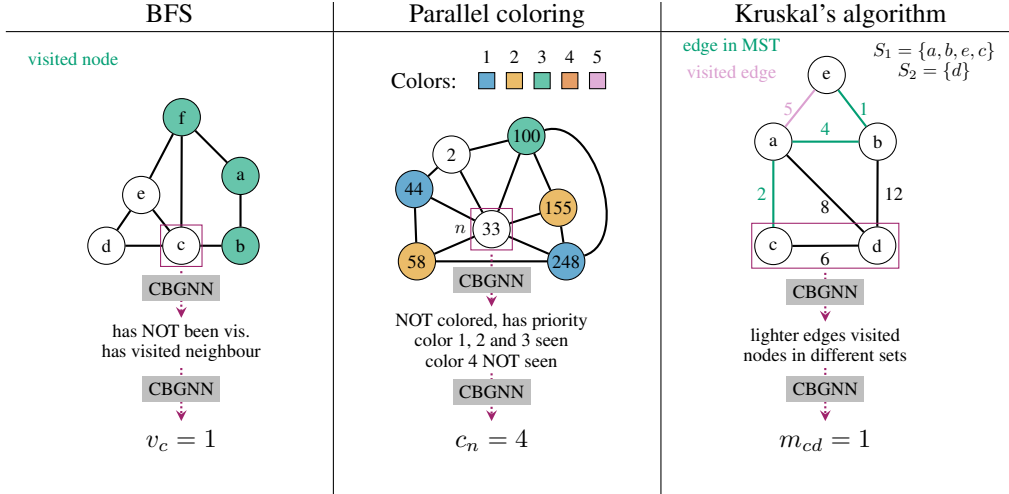


Figure 1: An overview of our Concept Bottleneck Graph Neural Network (CBGNN) approach. Importantly, CBGNN models can be trained to extract concept information for a given task as well as algorithm rules. We give examples of 3 algorithms, showing how CBGNN extract concepts from the input data and then uses these to compute the output.

provide model explanations in terms of human-understandable units, rather than individual features, pixels, or characters (e.g., the concepts of a wheel and a door are important for the detection of cars) (Kazhdan et al., 2020b). In particular, work on Concept Bottleneck Models (CBMs) relies on concepts and introduces a novel type of interpretable-by-design CNN, which perform input processing in two distinct steps: computing a set of concepts from an input, and then computing the output label from the concepts (Koh et al., 2020).

In this paper, we apply the idea of CBMs to GNN models, by introducing Concept Bottleneck Graph Neural Networks (CBGNNs). In particular, we rely on the *encode-process-decode* paradigm (Hamrick et al., 2018), and apply concept bottleneck layers before the output of GNN models – see Figure 1. By doing this we are able to extract update/termination rules for the step updates of *step-level* combinatorial optimisation approaches (Veličković et al., 2020b,a; Deac et al., 2020; Strathmann et al., 2021).

Importantly, we show that by relying on a suitable set of concepts and supervising on them, we are capable of deriving the rules of classical algorithms such as breadth-first search (Moore, 1959), and Kruskal’s algorithm (Kruskal, 1956), as well as more advanced heuristics such as parallel graph coloring (Jones and Plassmann, 1993). Furthermore, we present an approach to utilise node-level concepts for extracting graph-level rules. Our evaluation experiments demonstrate that all of our extracted rules strongly generalise to graphs of $5 \times$ larger size.

To summarise, we make the following contributions in this work:

- Present *Concept Bottleneck Graph Neural Networks* (CBGNN), a novel type of GNN relying on intermediate concept processing. To the best of our knowledge, this is the first work to apply concept bottleneck approaches to GNNs.
- Quantitatively evaluate our approach using three different case-studies (BFS, graph colouring, and Kruskal’s), showing that our CBGNN approach is capable of achieving performance on-par with that of existing state-of-the-art
- Qualitatively evaluate our approach, by demonstrating how the concepts utilised by CBGNN models can be used for providing rules summarising the heuristics the CBGNN has learned

2 Related work

GNN Explainability Recent work began exploring applications of XAI techniques in the context of GNNs. For instance, work in Pope et al. (2019); Baldassarre and Azizpour (2019); Schnake et al.

(2020) adapt feature-importance gradient-based approaches used for CNN applications (such as Class Activation Mappings, or Layer-wise Relevance Propagation) to GNNs, in order to identify the most important nodes/subgraphs responsible for individual predictions. Alternatively, works in Ying et al. (2019); Vu and Thai (2020); Luo et al. (2020) focus on more complex approaches unique to GNN explainability, such as those based on mutual information maximisation, or Markov blanket conditional probabilities of feature explanations. Importantly, these works focus on GNN tasks and benchmarks involving social networks, chemistry, or drug discovery, instead of focusing on combinatorial optimisation tasks, which is the focus of this work. Furthermore, these works focus on explaining pre-trained GNNs in a post-hoc fashion, whereas we focus on building GNN models interpretable-by-design. Finally, these works focus on feature-importance-based explanation approaches (i.e. returning relative importance of input nodes/subgraphs), whereas we rely on concept-based explanation approaches instead.

Concept-based Explainability A range of existing works have explored various concept-based explanations applied to CNN models. For instance, work in Ghorbani et al. (2019); Kazhdan et al. (2020b); Yeh et al. (2019) introduce approaches for extracting concepts from pre-trained CNNs in an unsupervised, or semi-supervised fashion. Work in Chen et al. (2020); Koh et al. (2020) rely on concepts for introducing CNN models interpretable-by-design, performing processing in two distinct steps: concept extraction, and label prediction. Other works on concepts include studying the connection between concepts and disentanglement learning (Kazhdan et al., 2021), as well as using concepts for data distribution shifts (Wijaya et al., 2021). Importantly, these works explore concepts exclusively in the context of CNNs, with Kazhdan et al. (2020a) being the only work exploring concepts in the context of RNN models. In this work, we focus on concept-based explainability for GNNs, where, similar to Koh et al. (2020), the concepts are human-specified.

Combinatorial Optimisation for GNNs Following the hierarchy defined in Cappart et al. (2021), our work classifies as a step-level approach. We directly extend on Veličković et al. (2020a,b), therefore we use the models presented in these works as baselines. We *do not* compare our model to an algorithm-level combinatorial optimisation approaches (Xu et al., 2020; Tang et al., 2020; Joshi et al., 2020) or unit-level ones (Yan et al., 2020) for the following reasons: *Algorithm-level approaches usually give one output per data sample* (rather than one output per step), but rules/invariants of a given algorithm come from how the iteration proceeds making algorithm-level combinatorial optimisation less suitable for a concept bottleneck. *Unit-level learning focuses on learning primitive units of computation*, such as taking maximum or merging lists and then combining these manually – having explanations at this level would not be of great benefit. To the best of our knowledge, only Veličković et al. (2020a) attempted to explain GNN predictions, using GNNExplainer (Ying et al., 2019). However, their model (i) was not explainable by design and (ii) required further optimisation *for a single sample* to give a *local explanation*. All other previous works operated in a black-box fashion and did not consider explainability of the learnt models.

3 Methodology

Encode-process-decode Following the “blueprint” for neural execution outlined in Veličković et al. (2020b), we model the algorithms by the encode-process-decode architecture (Hamrick et al., 2018). For each algorithm A , an *encoder network* f_A encodes the algorithm-specific node-level inputs $\mathbf{z}_i^{(t)}$ into the latent space. These node embeddings are then processed using the *processor network* P , usually a GNN. The processor takes as input the encoded inputs $\mathbf{Z}^{(t)} = \{\mathbf{z}_i^{(t)}\}_{i \in V}$ and graph edge index E to produce latent features $\mathbf{H}^{(t)} = \{\mathbf{h}_i^{(t)} \in \mathbb{R}^{|\mathcal{L}|}\}_{i \in V}$, where $|\mathcal{L}|$ is the size of the latent dimension. In contrast with previous work, we calculate algorithm outputs by first passing the latent embeddings through a *decoder network* g'_A , which produces concepts for each node $\mathbf{C}^{(t)} = \{\mathbf{c}_i^{(t)} \in (0, 1)^{|\mathcal{C}|}\}$, where $|\mathcal{C}|$ is number of concepts. The concepts are then passed through a *concept decoder* g_A to produce node-level outputs $\mathbf{Y}^{(t)} = \{\mathbf{y}_i^{(t)}\}$.

Where applicable, we also utilise a *termination network* T_A for deciding when to stop. However, in contrast with prior work, we observed that training is more stable if we calculate the termination probability based on potential *next step* embeddings (i.e. a belief about what is the state after an iteration has been executed). Additionally we found it insufficient to use the average node embeddings as input to T_A – averaging would obfuscate the signal if there is just a single node which should tell

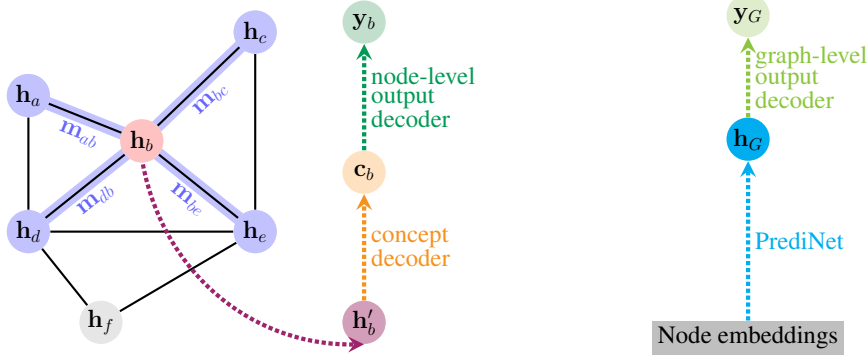


Figure 2: A high level overview of our GNN architecture: **Left:** (Eqns. 1-4) To produce node-level outputs, messages from neighbouring nodes (\mathbf{m}_{ij}) are combined with current node representation (\mathbf{h}_b), resulting in an updated representation (\mathbf{h}'_b). Concepts (\mathbf{c}_b) are then extracted from the updated representations, and node-level outputs (\mathbf{y}_b) are extracted from the concepts. **Right:** (Eqns. 5-8) A graph-level embedding (\mathbf{h}_G) is obtained by passing the node embeddings through PrediNet. We extract graph-level outputs \mathbf{y}_G (in our case – termination probability) directly from the latent state \mathbf{h}_G – graph-level concepts are extracted through a full enumeration approach over the node concepts.

us whether to continue iterating or not. Instead, we opted to use the output of an adapted PrediNet (Shanahan et al., 2020) architecture with one attention head. PrediNet is designed to represent the conjunction of elementary propositions, therefore it can (theoretically) capture the logical bias of the termination rules. The whole process is summarised in Figure 2 as well as in the equations below:

$$\begin{aligned}
 \mathbf{z}_i^{(t)} &= f_A(\mathbf{x}_i^{(t)}, \mathbf{h}_i^{(t-1)}) & (1) & \quad \mathbf{z}'_i^{(t)} = f_A(\mathbf{y}_i^{(t)}, \mathbf{h}_i^{(t)}) & (5) \\
 \mathbf{H}^{(t)} &= P(\mathbf{Z}^{(t)}, E) & (2) & \quad \mathbf{H}'^{(t)} = P(\mathbf{Z}'^{(t)}, E) & (6) \\
 \mathbf{c}_i^{(t)} &= \sigma(g'_A(\mathbf{z}_i^{(t)}, \mathbf{h}_i^{(t)})) & (3) & \quad \overline{\mathbf{H}}^{(t)} = \text{PrediNet}(\mathbf{H}'^{(t)}) & (7) \\
 \mathbf{y}_i^{(t)} &= g_A(\mathbf{c}_i^{(t)}) & (4) & \quad \tau^{(t)} = \sigma(T_A(\overline{\mathbf{H}}^{(t)})) & (8)
 \end{aligned}$$

where σ is a logistic sigmoid function. When using T_A , equations 1-8 are repeated if $\tau^{(t)} > 0.5$.

The combination of encoded inputs, together with the latent state of the given node, contains sufficient information not only about the output at a given step, but also: (i) a node’s current state and (ii) observations about other nodes’ states in its neighbourhood. If our concepts are engineered to capture some knowledge of either (i) or (ii), then we can extract meaningful algorithm output explanations *without providing any explicit information about how the algorithm works (theorems, invariants, etc.)*

Explicitly relational GNN architecture Some graph level tasks (e.g. deciding termination) can be reduced to a logical formula over all nodes – for the algorithms and concepts we consider, termination can be reduced to existence of a node with specific properties. (See graph-level rule extraction). We engineer this logical bias into the termination network τ by adapting PrediNet (Shanahan et al., 2020) to the graph domain. The PrediNet network architecture learns to represent conjunction/disjunction of elementary propositions and is therefore suitable for the termination task. We list the two minor modifications we made to adapt PrediNet to our tasks in Appendix A.

Extracting node-level algorithm rules Deciding node-level formulas for algorithm A is achieved by examining the weights of the concept decoder g_A . To achieve this, we used the open-source package `logic_explained_networks`¹ (Barbiero et al., 2021) implementing a wide collection of techniques to extract logic-based explanations from concept-bottleneck neural networks (Gori, 2017; Ciravegna et al., 2020). The library takes as inputs (i) node-level output decoder weights, (ii) predicted

¹Apache 2.0 Licence.

concepts from training data, and (iii) training data ground truth labels, and generates logic formulas in disjunctive normal form as outputs (Mendelson, 2009). By construction, the concept decoder $\bar{g}_A \approx g_A$ learnt from concepts \mathcal{C} to outputs \mathcal{O} is a Boolean map. As any Boolean function, it can be converted into a logic formula in disjunctive normal form by means of its truth-table (Mendelson, 2009). The weights of the concept decoder g_A are used to select the most relevant concepts for each output task. To get concise logic explanations *when many concepts are required as inputs* (in our experiments this is only the graph coloring task), we add a regularization term in the loss function minimising the $L1$ -norm of the concept decoder weights W , leading to sparse configurations of W . Later, at the training epoch t_{prune} , first-layer weights are pruned concept-wise, i.e. removing *all* the weights departing from the least relevant concepts:

$$\tilde{W}_j^1 = W_j^1 \mathbb{I}_{\|W_j^1\|_1 \geq \max_i \|W_i^1\|_1/2}, \quad \text{for } i = 1, \dots, |\mathcal{C}| \quad (9)$$

where \mathbb{I} is the indicator function and W^1 are weights of the first layer. Further details on logic extraction are provided in Appendix B.

Extracting algorithm termination rules When deciding whether to continue execution, we use the fact that a number of graph algorithms continue iterating until a node with a specific combination of concepts exists. Since it is unclear what combination of concepts we should supervise towards, we took a full enumeration approach when extracting rules for termination. First, we generate a sample j of the form (\mathcal{U}_j, τ'_j) from the training set from each step for a given graph. τ'_j is the ground truth for whether we should keep iterating, and $\mathcal{U}_j = \{\mathbf{c}'_1, \dots, \mathbf{c}'_k\}$ is a set of all unique concepts combinations,² *after the algorithm update state has been performed* (hence \mathbf{c}'). Given a set of concept indexes³ $\mathcal{I} \subseteq \mathcal{P}(\{1..|\mathcal{C}|\})$ and truth assignment $\mathcal{T} : \{1..|\mathcal{C}|\} \rightarrow \{0, 1\}$ telling us which concepts *must* be true/false, we check if the following is satisfied:

$$\forall j \left(\tau'_j = 1 \iff (\exists \mathbf{c} \in \mathcal{U}_j. \forall \mathcal{I}_i \in \mathcal{I}. \mathbf{c}_{\mathcal{I}_i} = \mathcal{T}(\mathcal{I}_i)) \right) \quad (10)$$

i.e. we should continue iterating if a special concept combination exists, and we should stop iterating if it does not. We employ a brute-force approach for finding \mathcal{I} and \mathcal{T} , breaking ties by preferring smaller \mathcal{I} .⁴ The complexity of such an approach is exponential, but if the concept bottleneck is carefully engineered, the number of necessary concepts, and/or the number of concepts in the special combination will be small, making the computation feasible.

More importantly, if the same enumeration approach was applied to the raw node input data a \mathcal{I}/\mathcal{T} combination *may not exist*. For example, the node-level inputs for the BFS task on each step do not tell us which nodes have visited neighbours (crucial for deciding termination). Additionally, if we have larger number of input features, the brute-force approach *may not be computationally feasible* – the combinations scale exponentially with the number of node features and concepts are one way to reduce this number.

4 Experimental setup

The code for our experiments can be found at <https://github.com/HekpoMaH/algorithmic-concepts-reasoning>.

Algorithms considered We apply our GNN to the following algorithms: *breadth-first search* (BFS), *parallel coloring* (Jones and Plassmann, 1993), a graph coloring heuristic, and *Kruskal’s minimum spanning tree* (MST) algorithm (Kruskal, 1956). BFS is modelled as a binary classification problem where we predict whether a node is visited or not, parallel coloring – as a classification task over the classes of possible node colors, plus one class for uncolored nodes. Kruskal’s is modelled as two tasks trained in parallel – one, as a classification task to choose the next edge to be considered for the MST and one to help us decide which nodes belong to the same set. As the original Kruskal’s algorithm executes for $|E|$ steps, *we do not learn termination for the MST task* and stick to a fixed number of steps. We show how we model inputs/outputs for all algorithms in Appendix C.

² k may vary across samples

³ $\{a..b\}$ denotes the set of integers from a to b

⁴In our case this broke all ties, but, if necessary, one can add tie-break on truth assignment, by e.g. choosing assignments with more true/false values

Table 1: Algorithms and their corresponding concepts. We provide some sample ground truth explanations. Visual examples of how the algorithms work can be seen in Figure 1.

Algorithm	Concepts	Example ground-truth explanations (not provided to the model)
BFS	<i>hasBeenVisited</i> (<i>hBV</i>) <i>hasVisitedNeighbours</i> (<i>hVN</i>)	$hVN(i) \implies y_i^{(t)} = 1$ $\exists i. \neg hBV(i) \wedge hVN(i) \implies \tau^{(t)} = 1$
Coloring	<i>isColored</i> (<i>iC</i>), <i>hasPriority</i> (<i>hP</i>) <i>colorXSeen</i> (<i>cXS</i>), $X \in \{1, \dots, 5\}$	$iC(i) \wedge c1S(i) \wedge \neg c2S(i) \implies y_i^{(t)} = 2$ $(\neg iC(i) \wedge hP(i) \wedge c1S(i) \wedge c2S(i) \wedge \neg c3S(i)) \implies y_i^{(t)} = 3$
Kruskal’s	<i>lighterEdgesVisited</i> (<i>lEV</i>) <i>nodesInSameSet</i> (<i>nISS</i>) <i>edgeInMst</i> (<i>eIM</i>)	$(lEV(i) \wedge \neg nISS(i) \wedge \neg eIM(i)) \implies y_i^{(t)} = 1$ $(nISS(i) \wedge \neg eIM(i)) \implies y_i^{(t)} = 0$

Importantly, all algorithms we experimented with possess the following two properties: (i) node/edge outputs are discrete and can be described in terms of concepts; (ii) continuing the execution can be reduced to the existence of a node with a specific combination of features. Examples of classical algorithms that *do not* fall into this category are the class of shortest path algorithms: to explain such algorithms, we would need to use arithmetic (e.g. minimum, sum) for the rules – something that concepts cannot directly capture. We leave explanation of such algorithms for future work.

To generate our concepts, we took into account what properties of the nodes/neighbourhood the algorithm uses, but we did not provide any details to the model *how* to use them. Table 1 gives more details on what concepts we chose and some example explanations. We give an outline how one can use these concepts for explaining the algorithms, in Appendix D.

Data generation Following prior research on the topic of neural execution (Veličković et al., 2020b), for BFS we generate graphs from a number of categories – ladder, grid, Erdős-Rényi (Erdős and Rényi, 1960), Barabási-Albert (Albert and Barabási, 2002), 4-Community graphs, 4-Caveman graphs and trees. For the coloring task, we limit the number of colors to 5 and then generate graphs where *all* nodes have fixed degree 5. This made the task both challenging (i.e. there are occasions where 5 colors are necessary) and feasible (we can generate graphs that are 5-colorable). Training data for these tasks graph size is fixed at 20 and we test with graph sizes of 20, 50 and 100 nodes. For Kruskal’s algorithm, we reused most graph categories for the BFS task, except the last three where the graph is either a tree or is not connected. Due to GPU memory constraints, training MST on graphs of size 20 required reducing the batch size by a factor of 4 and making the training very time consuming. Therefore, for the MST task we reduced the size of the training graphs to 8. Testing is still performed on graphs of size 20, 50 and 100. For all tasks we did a 10:1:1 train:validation:testing split.⁵ More details about the data generation are present in Appendix E.

Architectures tested We decided to choose message-passing neural networks (Gilmer et al., 2017) with the max aggregator for the main skeleton of our processor (GNN) architecture as this type of GNN is known to align well with algorithmic execution (Veličković et al., 2020b; Georgiev and Liò, 2020; Veličković et al., 2020a). However, due to the nonlinear nature of some of the tasks (parallel coloring) and the added concept bottleneck we found it beneficial to add a hidden layer to some of the encoders and decoders, rather than simply model them as an affine projection.

The Kruskal’s algorithm consists of several steps – masking out visited edges, finding the minimal edge from the unmasked and checking if two nodes are in the same set and unifying if they are not. The architecture for this algorithm, follows the main ideas of Figures 1&2, to implement them we

⁵When working with multiple graph categories, the ratio is preserved across each category

Table 2: Parallel coloring accuracies over 5 runs

Model	Metric	$ V = 20$	$ V = 50$	$ V = 100$
Standard	mean-step acc.	99.09 \pm 0.86%	98.74 \pm 0.44%	97.92 \pm 1.50%
	last-step acc.	99.25 \pm 0.56%	99.17 \pm 0.20%	99.13 \pm 0.29%
	term. acc.	98.79 \pm 0.86%	96.79 \pm 1.53%	95.08 \pm 2.89%
Bottleneck (+L1 and prune)	mean-step acc.	99.71 \pm 0.11%	99.23 \pm 0.21%	98.92 \pm 0.59%
	last-step acc.	99.69 \pm 0.13%	99.17 \pm 0.23%	99.10 \pm 0.22%
	term. acc.	99.61 \pm 0.18%	99.02 \pm 0.43%	98.59 \pm 0.77%
	formula mean-step acc.	99.71 \pm 0.12%	99.24 \pm 0.21%	98.93 \pm 0.59%
	formula last-step acc.	99.69 \pm 0.13%	99.16 \pm 0.22%	99.08 \pm 0.19%
	formula term. acc.	99.51 \pm 0.17%	99.02 \pm 0.43%	98.48 \pm 0.74%
	*concepts mean-step acc.	99.85 \pm 0.05%	99.60 \pm 0.10%	99.45 \pm 0.29%
	*concepts last-step acc.	99.72 \pm 0.07%	99.35 \pm 0.23%	99.42 \pm 0.09%

combine the architecture of Yan et al. (2020) for the first two steps and Veličković et al. (2020a) for the third step. More details can be found in Appendix F.

Experimental details We train our models using teacher forcing (Williams and Zipser, 1989) for a fixed number of epochs (500 for BFS, 3000 for the parallel coloring, 100 for Kruskal’s). When testing BFS/parallel coloring, we pick the model with the lowest sum of validation losses and when testing Kruskal’s – the model with highest last-step accuracy. For training we use Adam optimizer (Kingma and Ba, 2015) with initial learning rate of 0.001 and batch size 32. We optimise the sum of losses on the concept, output and termination (except for Kruskal’s, see above) predictions – for more details on how we define our losses see Appendix G. We evaluate the ability to strongly generalise on graphs with sizes 50 and 100. Standard deviations are obtained over 5 runs. For parallel coloring we add L1 regularisation and pruning on epoch 2000 to obtain higher quality explanations since every combination of (concepts, output) pair may not be observed during training. Libraries, code, and computing details are described in Appendix L. All hyperparameters were tuned manually.

Metrics We use a variety of metrics, such as *mean-step accuracy* (average accuracy of per-step outputs), *last-step accuracy* (average accuracy of final algorithm outputs) and *termination accuracy* (average accuracy of predicting termination). Similarly, we define: *concepts mean-step accuracy* and *concepts last-step accuracy* as well *formula mean-step accuracy*, *formula last-step accuracy* and *formula termination accuracy*. The last three are derived by applying the *extracted formulas* to the *predicted concepts* for predicting the output/termination instead of using the respective neural network. The motivation behind is that if we achieve high concept accuracies and high formula accuracies then the formulas are likely to be representing the underlying algorithm (or data) accurately.

Qualitative analysis We provide several qualitative experiments: (i) We fit a decision tree (DT) for the $\mathcal{C} \rightarrow \mathcal{O}$ task ($\mathcal{C} \rightarrow \mathcal{T}$ is not possible, due to DTs working on fixed size node-level features). Concepts and targets are obtained from the *ground truth* concepts and target classes of *all training data nodes* at each step for each graph. (ii) We also plot the concepts last/mean step accuracy vs epoch for each concept and provide further analysis on which concept the networks find the most difficult. (iii) We provide sample target class explanations for each algorithm.

5 Results and discussion

Concept accuracies As can be seen from Tables 2&3 and Table 6, Appendix H (metrics with an asterisk) we are able to learn concepts with high accuracy (99% and higher accuracy for BFS and parallel coloring). Results show that GNNs are capable of producing high-level concepts, capturing either node or neighbourhood information, for these algorithmic tasks and the learned concept extractors strongly generalise – concept accuracy does not drop even for $5\times$ larger graphs.

Parallel algorithms: BFS and coloring For the BFS task both the baseline and bottlenecked model perform optimally in line with the state of the art. We therefore present results from the BFS

Table 3: Kruskal’s algorithm accuracies over 5 runs

Model	Metric	$ V = 20$	$ V = 50$	$ V = 100$
Standard	mean-step acc.	96.75 \pm 0.15%	95.41 \pm 0.09%	94.68 \pm 0.10%
	last-step acc.	93.70 \pm 0.33%	90.10 \pm 2.80%	86.69 \pm 4.28%
Bottleneck	mean-step acc.	96.93 \pm 0.13%	95.86 \pm 0.37%	95.27 \pm 0.59%
	last-step acc.	94.00 \pm 0.24%	92.20 \pm 0.52%	91.29 \pm 0.86%
	formula mean-step acc.	96.79 \pm 0.37%	95.77 \pm 0.31%	95.25 \pm 0.54%
	formula last-step acc.	93.70 \pm 0.71%	91.92 \pm 0.47%	91.15 \pm 0.60%
	*concepts mean-step acc.	97.91 \pm 0.08%	97.21 \pm 0.22%	96.80 \pm 0.35%
	*concepts last-step acc.	99.56 \pm 0.29%	99.49 \pm 0.49%	97.09 \pm 0.29%

Table 4: Sample explanations for each algorithm *obtained from the learned model*. *cXS* denotes *color X Seen*, *nISS* is *nodesInSameSet*, *LEV* is *lighterEdgesVisited*, *eIM* is *edgeInMst*.

Algorithm	Thing to explain	Explanation
BFS	n is visited	$hasVisitedNeighbours(n)$
	continue execution	$\exists n. \neg hasBeenVisited(n) \wedge hasVisitedNeighbours(n)$
parallel coloring	n has color 2	$(isColored(n) \wedge \neg hasPriority(n) \wedge c1S(n) \wedge \neg c2S(n)) \vee (hasPriority(n) \wedge c1S(n) \wedge \neg c2S(n) \wedge \neg isColored(n))$
	n has color 5	$(isColored(n) \wedge \neg hasPriority(n) \wedge c1S(n) \wedge c2S(n) \wedge c3S(n) \wedge c4S(n)) \vee (hasPriority(n) \wedge c1S(n) \wedge c2S(n) \wedge c3S(n) \wedge c4S(n) \wedge \neg isColored(n))$
	continue execution	$\exists n. \neg isColored(n)$
Kruskal’s	e not in MST	$(nISS(e) \wedge \neg eIM(e)) \vee (\neg LEV(e) \wedge \neg eIM(e))$
	e in MST	$(LEV(e) \wedge nISS(e) \wedge eIM(e)) \vee (LEV(e) \wedge \neg nISS(e) \wedge \neg eIM(e))$

task in Appendix H. Results from the parallel coloring task are shown in Table 2. Apart from the high accuracy achieved, our results show that: (i) the bottleneck doesn’t have a major impact on the final model accuracy – original metrics⁶ remain the same or are better for both algorithms; (ii) we are able to learn concepts accurately and (iii) the extracted rules are accurate – applying them to the accurately predicted concepts in order to produce output has no significant negative effect on the predictive accuracy of our model – formula based accuracies do not deviate more than 5-6% than the original metrics.

Qualitative analysis: decision trees We visualise the fitted decision trees (DTs) for each algorithm in Appendix I. In all cases the logic of the DT follows the logic of the original algorithm. Additionally, the leaf nodes of all decision trees contain samples from a single class showing that concepts were capable of capturing the complexity of the algorithm.

Qualitative analysis: concept learning curves We present per concept learning curves for the parallel coloring in Figure 3 and for Kruskal’s in Figure 4: (i) Parallel coloring exhibits many occasions where there are drops of concept accuracy across almost all concepts. If we observe more carefully Figure 3a, we will notice that they coincide with a drop of the accuracy of *hasPriority* concept. This drop also explains the lower last-step concept accuracy – changing the coloring order early on may produce quite different final coloring. To confirm this observations, we trained an oracle model that is always provided with the correct value for *hasPriority*. Such oracle model achieved almost perfect concept accuracy – we provide a plot of the concept learning curves in Appendix J; (ii) The concept instability was present only in the beginning for Kruskal’s, but it converged to a stable solution. The reason *edgeInMst* concept remained with the lowest last-step accuracy is that the overall last-step accuracy of the model is lower.

⁶namely mean-, last-step accuracy and termination accuracy

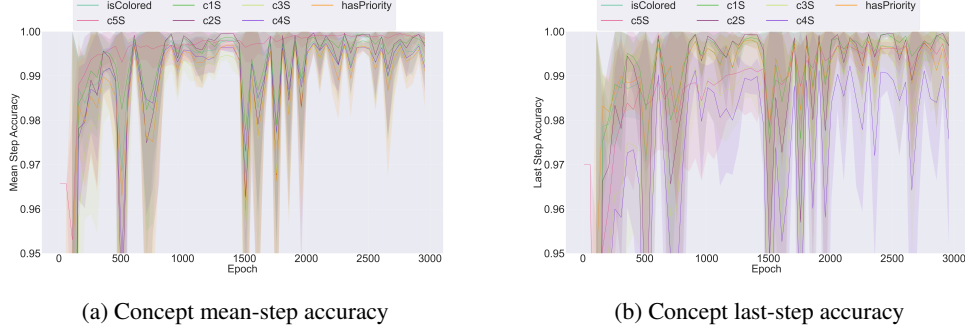


Figure 3: Concept accuracies per epoch of the parallel coloring algorithm. (1 point every 50 epochs). cXS is *colorXSeen*. Note the y axis scale. It can be observed that the *hasPriority* concept is one of the worst performing concepts. This leads to nodes being colored in a different order and therefore lower last-step concept accuracy for concepts related to colors. Standard deviation obtained from 5 runs.

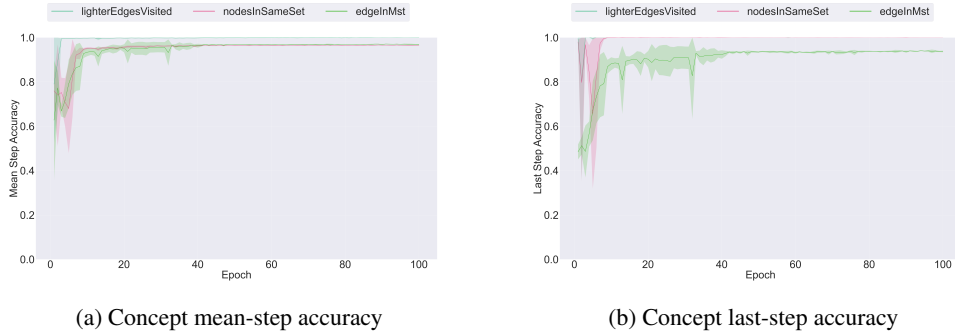


Figure 4: Concept accuracies per epoch of the Kruskal's algorithm on graphs with 20 nodes. (1 point per epoch). After an initial instability concepts are consistently accurate. Standard deviation obtained from 5 runs.

Qualitative analysis: explanations We list examples of obtained explanations in Table 4 and present *all* explanations obtained from the algorithms in in Appendix K. The extracted rules show that concepts are one way to extract accurate representation of the rules of the algorithm. E.g. we can (re)infer from the listed rules for parallel coloring that for getting a given color that color should not be seen in the neighbourhood and colors coming before that one have already been seen.

We additionally observed, that as the number of concepts increases, if we need shorter and more general rules we need more and more data. One way to alleviate such problem is L1 regularisation and pruning – we additionally perform an ablation study in Appendix K showing that without regularisation rules are still usable (giving good formula accuracy) *but are less general*.

6 Conclusions

We presented concept-based reasoning on graph algorithms through Concept Bottleneck Graph Neural Networks. We demonstrated through the surveyed algorithms, that we can accurately learn node-level concepts without impacting performance. Moreover, by examining training data and model weights, we are capable of explaining each node-level output classes with formulas based on the defined concepts. Concepts also allow us perform a *unsupervised* rule extraction of certain graph-level tasks, such as deciding when to terminate. Extracted rules are interpretable and applying them does not heavily impact accuracy.

References

- Adadi, A. and Berrada, M. (2018). Peeking inside the black-box: a survey on explainable artificial intelligence (xai). *IEEE access*, 6:52138–52160.
- Albert, R. and Barabási, A.-L. (2002). Statistical mechanics of complex networks. *Reviews of modern physics*, 74(1):47.
- Baldassarre, F. and Azizpour, H. (2019). Explainability techniques for graph convolutional networks. *arXiv preprint arXiv:1905.13686*.
- Barbiero, P., Ciravegna, G., Georgiev, D., and Giannini, F. (2021). Lens: a python library for logic explained networks. *arXiv preprint*.
- Battaglia, P. W., Pascanu, R., Lai, M., Rezende, D. J., and Kavukcuoglu, K. (2016). Interaction networks for learning about objects, relations and physics. In Lee, D. D., Sugiyama, M., von Luxburg, U., Guyon, I., and Garnett, R., editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 4502–4510.
- Cappart, Q., Chételat, D., Khalil, E., Lodi, A., Morris, C., and Veličković, P. (2021). Combinatorial optimization and reasoning with graph neural networks. *arXiv preprint arXiv:2102.09544*.
- Chen, Z., Bei, Y., and Rudin, C. (2020). Concept whitening for interpretable image recognition. *Nature Machine Intelligence*, 2(12):772–782.
- Ciravegna, G., Giannini, F., Melacci, S., Maggini, M., and Gori, M. (2020). A constraint-based approach to learning and explanation. In *AAAI*, pages 3658–3665.
- Deac, A., Velickovic, P., Milinkovic, O., Bacon, P., Tang, J., and Nikolic, M. (2020). XLVIN: executed latent value iteration nets. *CoRR*, abs/2010.13146.
- Doshi-Velez, F. and Kim, B. (2017). Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*.
- Erdős, P. and Rényi, A. (1960). On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci*, 5(1):17–60.
- Galler, B. A. and Fisher, M. J. (1964). An improved equivalence algorithm. *Communications of the ACM*, 7(5):301–303.
- Georgiev, D. and Liò, P. (2020). Neural bipartite matching. In *Graph Representation Learning and Beyond (GRL+) workshop*.
- Ghorbani, A., Wexler, J., Zou, J., and Kim, B. (2019). Towards automatic concept-based explanations. *arXiv preprint arXiv:1902.03129*.
- Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. (2017). Neural message passing for quantum chemistry. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 1263–1272.
- Gori, M. (2017). *Machine Learning: A constraint-based approach*. Morgan Kaufmann.
- Hamrick, J. B., Allen, K. R., Bapst, V., Zhu, T., McKee, K. R., Tenenbaum, J., and Battaglia, P. W. (2018). Relational inductive bias for physical construction in humans and machines. In *Proceedings of the 40th Annual Meeting of the Cognitive Science Society, CogSci 2018, Madison, WI, USA, July 25-28, 2018*.
- Jones, M. and Plassmann, P. (1993). A parallel graph coloring heuristic. *SIAM J. Sci. Comput.*, 14:654–669.
- Joshi, C. K., Cappart, Q., Rousseau, L., Laurent, T., and Bresson, X. (2020). Learning TSP requires rethinking generalization. *CoRR*, abs/2006.07054.
- Kazhdan, D., Dimanov, B., Jamnik, M., and Liò, P. (2020a). Meme: Generating rnn model explanations via model extraction. *arXiv preprint arXiv:2012.06954*.

- Kazhdan, D., Dimanov, B., Jamnik, M., Liò, P., and Weller, A. (2020b). Now you see me (cme): Concept-based model extraction. *arXiv preprint arXiv:2010.13233*.
- Kazhdan, D., Dimanov, B., Terre, H. A., Jamnik, M., Liò, P., and Weller, A. (2021). Is disentanglement all you need? comparing concept-based & disentanglement approaches. *arXiv preprint arXiv:2104.06917*.
- Kingma, D. P. and Ba, J. (2015). Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Koh, P. W., Nguyen, T., Tang, Y. S., Musmann, S., Pierson, E., Kim, B., and Liang, P. (2020). Concept bottleneck models. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 5338–5348. PMLR.
- Kruskal, J. B. (1956). On the shortest spanning subtree of a graph and the traveling salesman problem. *Proceedings of the American Mathematical society*, 7(1):48–50.
- Luo, D., Cheng, W., Xu, D., Yu, W., Zong, B., Chen, H., and Zhang, X. (2020). Parameterized explainer for graph neural network. *arXiv preprint arXiv:2011.04573*.
- McCluskey, E. J. (1956). Minimization of boolean functions. *The Bell System Technical Journal*, 35(6):1417–1444.
- McColl, H. (1878). The calculus of equivalent statements (third paper). *Proceedings of the London Mathematical Society*, 1(1):16–28.
- Mendelson, E. (2009). *Introduction to mathematical logic*. CRC press.
- Molnar, C. (2020). *Interpretable machine learning*. Lulu. com.
- Moore, E. F. (1959). The shortest path through a maze. In *Proc. Int. Symp. Switching Theory, 1959*, pages 285–292.
- Pal, A., Eksombatchai, C., Zhou, Y., Zhao, B., Rosenberg, C., and Leskovec, J. (2020). Pinnersage: Multi-modal user embedding framework for recommendations at pinterest. In Gupta, R., Liu, Y., Tang, J., and Prakash, B. A., editors, *KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020*, pages 2311–2320. ACM.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al. (2019). Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint arXiv:1912.01703*.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., et al. (2011). Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830.
- Pope, P. E., Kolouri, S., Rostami, M., Martin, C. E., and Hoffmann, H. (2019). Explainability methods for graph convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10772–10781.
- Prüfer, H. (1918). Neuer beweis eines satzes über permutationen. *Arch. Math. Phys*, 27(1918):742–744.
- Quine, W. V. (1952). The problem of simplifying truth functions. *The American mathematical monthly*, 59(8):521–531.
- Schnake, T., Eberle, O., Lederer, J., Nakajima, S., Schütt, K., Müller, K., and Montavon, G. (2020). Higher-order explanations of graph neural networks via relevant walks. *arXiv: 2006.03589*.
- Shanahan, M., Nikiforou, K., Creswell, A., Kaplanis, C., Barrett, D. G. T., and Garnelo, M. (2020). An explicitly relational neural network architecture. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 8593–8603. PMLR.

- Stokes, J. M., Yang, K., Swanson, K., Jin, W., Cubillos-Ruiz, A., Donghia, N. M., MacNair, C. R., French, S., Carfrae, L. A., Bloom-Ackermann, Z., Tran, V. M., Chiappino-Pepe, A., Badran, A. H., Andrews, I. W., Chory, E. J., Church, G. M., Brown, E. D., Jaakkola, T. S., Barzilay, R., and Collins, J. J. (2020). A Deep Learning Approach to Antibiotic Discovery. *Cell*, 180(4):688–702.e13.
- Strathmann, H., Barekatin, M., Blundell, C., and Veličković, P. (2021). Persistent message passing. In *ICLR 2021 Workshop on Geometrical and Topological Representation Learning*.
- Tang, H., Huang, Z., Gu, J., Lu, B.-L., and Su, H. (2020). Towards scale-invariant graph-related problem solving by iterative homogeneous gnns. *Advances in Neural Information Processing Systems*, 33.
- Veličković, P., Buesing, L., Overlan, M. C., Pascanu, R., Vinyals, O., and Blundell, C. (2020a). Pointer graph networks. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Veličković, P., Ying, R., Padovano, M., Hadsell, R., and Blundell, C. (2020b). Neural execution of graph algorithms. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Vu, M. N. and Thai, M. T. (2020). Pgm-explainer: Probabilistic graphical model explanations for graph neural networks. *arXiv preprint arXiv:2010.05788*.
- Watts, D. J. (1999). Networks, dynamics, and the small-world phenomenon. *American Journal of sociology*, 105(2):493–527.
- Wijaya, M. A., Kazhdan, D., Dimanov, B., and Jamnik, M. (2021). Failing conceptually: Concept-based explanations of dataset shift. *arXiv preprint arXiv:2104.08952*.
- Williams, R. J. and Zipser, D. (1989). A learning algorithm for continually running fully recurrent neural networks. *Neural computation*, 1(2):270–280.
- Xu, K., Li, J., Zhang, M., Du, S. S., Kwarabazashi, K., and Jegelka, S. (2020). What can neural networks reason about? In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Yan, Y., Swersky, K., Koutra, D., Ranganathan, P., and Hashemi, M. (2020). Neural execution engines: Learning to execute subroutines. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H., editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Yeh, C.-K., Kim, B., Arik, S. O., Li, C.-L., Pfister, T., and Ravikumar, P. (2019). On completeness-aware concept-based explanations in deep neural networks. *arXiv preprint arXiv:1910.07969*.
- Ying, R., Bourgeois, D., You, J., Zitnik, M., and Leskovec, J. (2019). Gnnexplainer: Generating explanations for graph neural networks. *Advances in neural information processing systems*, 32:9240.

A Adapting PrediNet to graphs

As already discussed, PrediNet (Shanahan et al., 2020) is an architecture that operates on the image and learns to capture logical bias. To adapt it to the graph domain and more specifically the domain of our tasks, we made the following two modifications: (i) Instead of flattening a convoluted image we use global graph pooling (max pooling) for the calculation of the query and (ii) we do not concatenate any positional information to the input feature vectors. Everything else is kept the same as the original PrediNet:

$$\mathbf{Q}_1^h = \text{pool}(\mathbf{H}'^{(t)})W_{Q_1}^h \quad (11)$$

$$\mathbf{Q}_2^h = \text{pool}(\mathbf{H}'^{(t)})W_{Q_2}^h \quad (12)$$

$$\mathbf{K} = \mathbf{H}'^{(t)}W_k \quad (13)$$

$$\mathbf{E}_1^h = \text{softmax}(\mathbf{Q}_1^h K^\top) \mathbf{H}'^{(t)} \quad (14)$$

$$\mathbf{E}_2^h = \text{softmax}(\mathbf{Q}_2^h K^\top) \mathbf{H}'^{(t)} \quad (15)$$

$$\mathbf{D}^h = \mathbf{E}_1^h W_S - \mathbf{E}_2^h W_S \quad (16)$$

$$\text{out} = \text{Linear}(\mathbf{D}^h) \quad (17)$$

B Extraction of formulas

By construction, the function $\bar{f} \approx f$ learnt from concepts \mathcal{C} to outputs \mathcal{O} is a Boolean map. As any Boolean function, the map $\bar{f} : \mathcal{C} \rightarrow \mathcal{O}$ can be converted into a First-Order Logic (FOL) formula φ in Disjunctive Normal Form (DNF) (Mendelson, 2009) by means of its truth-table. The truth table provides a formal mechanism to extract logic rules of increasing complexity for individual observations, for cluster of samples, or for a whole class.

FOL extraction A logic formula can be extracted for each sample $c \in \mathcal{C}$ corresponding to a single row of the truth table:

$$\varphi_c = \tilde{c}_1 \wedge \dots \wedge \tilde{c}_k \quad \text{where } \tilde{c}_j := \begin{cases} \bar{c}_j, & \text{if } c_j \geq t \\ \neg \bar{c}_j, & \text{if } c_j < t \end{cases}, \text{ for } j = 1, \dots, k, \quad (18)$$

Example-level formulas can be formally aggregated providing logic rules φ_S for a cluster of samples $S \subseteq S^* = \{\bar{c} \mid \bar{f}(\bar{c}) = t\}$ of the same class t

$$\varphi_S = \bigvee_{c \in S} \varphi_c = \bigvee_{c \in S} \tilde{c}_1 \wedge \dots \wedge \tilde{c}_k \quad (19)$$

Any formula φ_S can be thought of as a $\{0, 1\}$ -valued mapping defined on \mathcal{C} , that is equal to 1 exactly on a cluster S . We may also get an explicit explanation relating φ_S and the Boolean function \bar{f} by the FOL formula:

$$\forall \bar{c} \in S : \varphi_S(\bar{c}) \rightarrow \bar{f}(\bar{c})$$

A formula for a whole class t (e.g. $t = 1$) can be formally obtained using Eq. 19 by considering S^* , i.e. aggregating all the example-level rules for the same class t . Also in this case it is possible to get a FOL formula relating φ and the Boolean map \bar{f} by means of:

$$\forall \bar{c} \in S^* : \bar{f}(\bar{c}) \leftrightarrow \varphi(\bar{c})$$

FOL simplification The aggregation of many example-level explanations may increase the length and complexity of the FOL formula being extracted for a cluster of samples or for a whole class. However, existing techniques as the Quine–McCluskey algorithm can be used to get compact and simplified equivalent FOL expressions (McColl, 1878; Quine, 1952; McCluskey, 1956). For instance, the explanation " $(\text{person} \wedge \text{nose}) \vee (\neg \text{person} \wedge \text{nose})$ " can be formally simplified in " nose ".

C Modelling algorithmic reasoning

We begin the modelling by designing the initial states for the algorithm:

$$\text{BFS} : x_i^{(1)} = \begin{cases} 1 & i = s \\ 1 & i \neq s \\ 0 & \text{otherwise} \end{cases} \quad \text{coloring} : x_i^{(1)} = (0, \text{binary}(p_i)) \quad \text{Kruskal's} : x_i^{(1)} = (0, w_i) \quad (20)$$

where s is a randomly chosen starting node, p_i denotes a randomly chosen⁷ *constant* integer priority for node i and w_i denotes the weight of edge i . For parallel coloring we leave class 0 for denoting uncolored nodes in the coloring task. The extra 0 in the input for Kruskal's will be a bit value whether an edge is selected to be in the minimum spanning tree. Similar to Yan et al. (2020) we represent numbers in binary, and learn an embedding for each bit position. We also one hot encode any class information, such as visited/unvisited or color of a node (or no color present).

Next step inputs are calculated according to the specific algorithm: for BFS, a node becomes reachable if it had visited neighbours⁸. For parallel coloring let us assume that there is a total ordering between colors, e.g. color 1 comes before color 2, and so on. A node keeps its color if it was previously colored and becomes colored with the first unseen color in the neighbourhood if it has the highest priority of all uncolored nodes in the neighbourhood. For Kruskal's algorithm an edge remains in the MST it has been already selected and an edge is selected to be inserted in the MST if all lighter edges have been checked and the nodes, connected by that edge are not in the same set. Note that for Kruskal's $\mathbf{x}^{(t)}$ ranges over the edges:

$$\text{BFS} : x_i^{(t+1)} = \begin{cases} 1 & x_i^{(t)} = 1 \\ 1 & \exists j. (j, i) \in E \wedge x_j^{(t)} = 1 \\ 0 & \text{otherwise} \end{cases} \quad (21)$$

$$\text{coloring} : x_i^{(t+1)} = \begin{cases} x_i^{(t)} & \text{col}_i^{(t)} \neq 0 \\ (\text{col}_i^{(t+1)}, p_i) & x_i^{(t)} = 0 \wedge p_i = \max\{p_j, j \in \mathcal{N}'_i \cup \{i\}\} \\ & \wedge \text{col}_i^{(t+1)} = \min\{\text{col}, \text{col} \notin \text{seen}_j\} \\ (0, p_i) & \text{otherwise} \end{cases} \quad (22)$$

$$\text{Kruskal's} : x_i^{(t+1)} = \begin{cases} x_i^{(t)} & \text{if already selected} \\ (1, w_i) & \text{if lighter edges visited} \wedge (a, b) = \text{edge}_i \\ & \text{set}^{(t)}(a) \neq \text{set}^{(t)}(b) \\ (0, w_i) & \text{otherwise} \end{cases} \quad (23)$$

where for the coloring task $\text{col}_i^{(t)}$ denotes the coloring of node i at step (t) , \mathcal{N}'_i denotes the set of uncolored neighbours, and seen_i denotes the set of colors observed in the neighbourhood of node i . For Kruskal's, $(a, b) = \text{edge}_i$ denotes that edge i connects nodes a and b , and $\text{set}^{(t)}(n)$ denotes the set of node n at timestep t . Before we proceed to the next step, if edge $\text{edge}_i = (a, b)$ is added to the MST, we unify the sets of the two edge nodes, i.e. $\text{set}^{(t+1)}(a) = \text{set}^{(t+1)}(b) = \text{set}^{(t)}(a) \cup \text{set}^{(t)}(b)$.

The ground-truth outputs we supervise towards are $\hat{y}_i^{(t)} = x_i^{(t+1)}$ for BFS, $\hat{y}_i^{(t)} = \text{col}_i^{(t+1)}$ for the coloring task and $\hat{y}_i^{(t)} = x_i^{(t+1)}[0]$ for Kruskal's (i.e. 0/1 value whether each edge is in the MST). We do not aim to reconstruct input node priorities or input edge weights.

D Algorithms and concepts

Given the concepts algorithms in Table 5, here how each algorithm can be executed with its own concepts:

⁷We ensure no two nodes connected by an edge have the same priority

⁸As we utilise self-loops for retention of self-information, this covers the corner case of a starting node being visited without having visited neighbours

Table 5: Algorithms and their corresponding concepts. Explanations are added to give a notion of how concepts can be used.

Algorithm	Concepts	Example ground-truth explanations (not provided to the model)
BFS	<i>hasBeenVisited</i> (<i>hBV</i>) <i>hasVisitedNeighbours</i> (<i>hVN</i>)	$hVN(i) \implies y_i^{(t)} = 1$ $\exists i. \neg hBV(i) \wedge hVN(i) \implies \tau^{(t)} = 1$
coloring	<i>isColored</i> (<i>iC</i>), <i>hasPriority</i> (<i>hP</i>) <i>colorXSeen</i> (<i>cXS</i>), $X \in \{1, \dots, 5\}$	$iC(i) \wedge c1S(i) \wedge \neg c2S(i) \implies y_i^{(t)} = 2$ $(\neg iC(i) \wedge hP(i) \wedge c1S(i) \wedge c2S(i) \wedge \neg c3S(i)) \implies y_i^{(t)} = 3$
Kruskal's	<i>lighterEdgesVisited</i> (<i>lEV</i>) <i>nodesInSameSet</i> (<i>nISS</i>) <i>edgeInMst</i> (<i>eIM</i>)	$(lEV(i) \wedge \neg nISS(i) \wedge \neg eIM(i)) \implies y_i^{(t)} = 1$ $(lEV(i) \wedge nISS(i) \wedge \neg eIM(i)) \implies y_i^{(t)} = 0$

- BFS – a node is visited if it has a visited neighbour and the execution continues until there is an unvisited node with visited neighbours.
- parallel coloring – check if node has been colored, if not, check if it has priority to be colored on this step, if yes, check colors in the neighbourhood. Continue execution until a node with $\neg isColored$ exists.
- Kruskal's algorithm – check if lighter edges have been visited by the algorithm, if not, edge is not in MST. Otherwise, check if the nodes connected by this edge belong in the same set and then if the edge has been previously selected. If the nodes are in the same set, the edge is in the MST for the next iteration if it has been selected for the MST previously. If nodes are not in the same set (this implies that the edge has not been selected for the MST), then this edge is in the MST from now on.

E Data generation

In general we aimed to maintain 10:1:1 train:validation:test split ratio and to preserve the ratio across different data generation techniques.

For the BFS task we generate graphs from the following categories:

- *Ladder graphs
- *2D grid graphs – during generation we aimed the grid to be as close to square as possible
- Trees – uniformly generated from the Prüfer sequence (Prüfer, 1918)
- *Erdős-Rényi (Erdős and Rényi, 1960), p of edge is $\min\left(\frac{\log_2 |V|}{|V|}, 0.5\right)$
- *Barabási-Albert (Albert and Barabási, 2002) graphs, with either 4 or 5 edges attached to each incoming node
- 4-Community graphs – generated by creating 4 disjoint Erdős-Rényi graphs with edge probability 0.7 and then interconnecting their nodes with probability of 0.01
- 4-Caveman (Watts, 1999) – intra-clique edges are removed with $p = 0.7$ and $0.025|V|$ shortcut edges are inserted after that between cliques.

From each category we generate 100 graphs for training and 10 for validation/testing. In total that is 700 for training and 70 for validation/testing. To test strong generalisation, for a given number of nodes in the graph, we generate 10 more test graphs of each category.

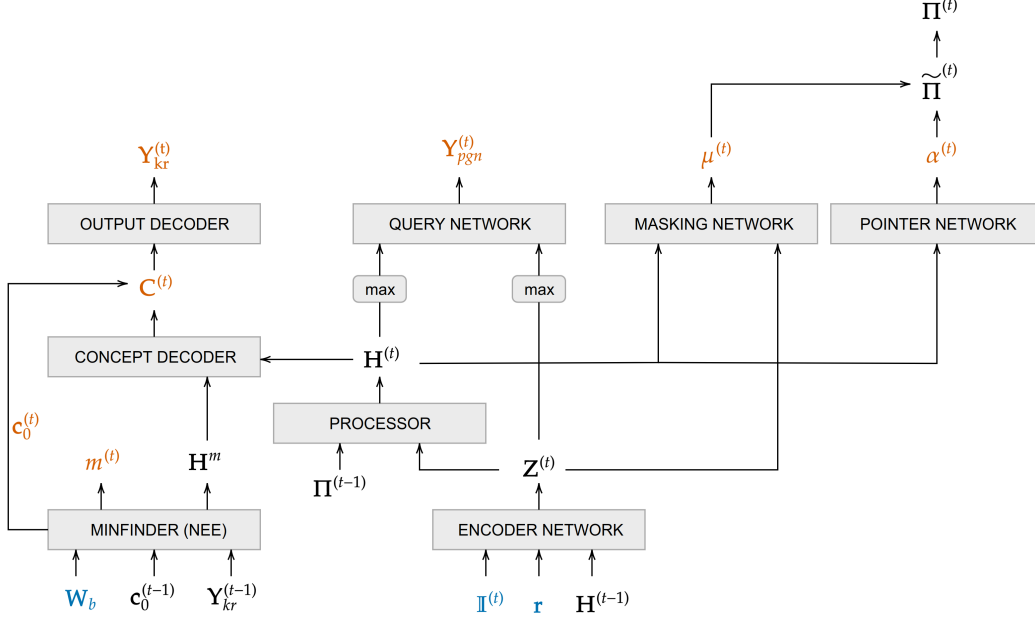


Figure 5: The dataflow for the Kruskal’s algorithm. Blue variables are input variables, black inputs come from previous steps, supervision is performed on orange outputs.

For the parallel coloring task we limit the number of colors to 5. We generate 800/80/80 graphs for training/validation/testing with nodes of fixed degree 5, so as to limit the possibility of 5 colors not being enough. We pick the priority p_i uniformly from the interval $[0, 255]$.

For Kruskal’s, we generate 500 graphs of size 8 of each category with a star from the list above for training and 50 for testing/validation of sizes $[8, 20, 50, 100]$.

To facilitate easier retention of self-information in the GNN, for the BFS and parallel coloring tasks, we insert self-loops to all nodes for a given graph similar to Veličković et al. (2020b). For Kruskal’s, all nodes start with a pointer to itself (see Appendix F)

F Implementing Kruskal’s algorithm

For the implementation of the Kruskal’s algorithm we need to implement a min finding subroutine with masking and a disjoint-set union (DSU) (Galler and Fisher, 1964) data-structure. We architecture the first as a Neural Execution Engine (NEE) (Yan et al., 2020) that takes in the binary edge weights \mathbf{W}_b , last step *lighterEdgesVisited* concept for all edges and last-step outputs $\mathbf{Y}_{kr}^{(t-1)}$. The NEE outputs a pointer to the next minimal edge $m^{(t)}$ and the next *lighterEdgesVisited* for all edges. We implement the DSU as a Pointer Graph Network (PGN) (Veličković et al., 2020a) that takes as input: (i) an indicator, which is the two nodes of the minimal edge on the current step (teacher-forced during training), (ii) node priorities (used as a tie-break when unifying sets) and (iii) what’s the last hidden state of the PGN algorithm $\mathbf{H}^{(t-1)}$. The PGN’s processor also takes the last step pointers $\Pi^{(t-1)}$ of the DSU data structure. The PGN predicts: (i) whether the two nodes are in the same set $\mathbf{Y}_{pgn}^{(t)}$; (ii) a mask $\mu^{(t)}$ over which nodes should change their DSU pointers and (iii) an estimate of the new pointer matrix. The next step pointers are updated from (ii) and (iii) as follows:

$$\tilde{\Pi}_{ij}^{(t)} = \mu_i^{(t)} \tilde{\Pi}_{ij}^{(t-1)} + (1 - \mu_i^{(t)}) \mathbb{I}_{j=\arg\max_k (\alpha_{ik}^{(t)})} \quad \Pi_{ij}^{(t)} = \tilde{\Pi}_{ij}^{(t)} \vee \tilde{\Pi}_{ji}^{(t)} \quad (24)$$

To predict the rest of the concepts⁹ for every edge $e = (u, v)$ we first concatenate the hidden state of the first NEE transformer \mathbf{H}_e^m with the hidden states of the PGN for nodes u and v , $\mathbf{H}_u^{(t)}$ and $\mathbf{H}_v^{(t)}$.

⁹apart from *lighterEdgesVisited*

Table 6: BFS accuracies over 5 runs

Model	Metric	$ V = 20$	$ V = 50$	$ V = 100$
Standard	mean-step acc.	100.0 \pm 0.0%	100.0 \pm 0.0%	100.0 \pm 0.0%
	last-step acc.	100.0 \pm 0.0%	100.0 \pm 0.0%	100.0 \pm 0.0%
	term. acc.	100.0 \pm 0.0%	100.0 \pm 0.0%	100.0 \pm 0.0%
Bottleneck (+next step pool)	mean-step acc	100.0 \pm 0.0%	100.0 \pm 0.0%	100.0 \pm 0.0%
	last-step acc	100.0 \pm 0.0%	100.0 \pm 0.0%	100.0 \pm 0.0%
	term. acc.	100.0 \pm 0.0%	100.0 \pm 0.0%	100.0 \pm 0.0%
	formula mean-step acc.	100.0 \pm 0.0%	100.0 \pm 0.0%	100.0 \pm 0.0%
	formula last-step acc.	100.0 \pm 0.0%	100.0 \pm 0.0%	100.0 \pm 0.0%
	formula term. acc.	100.0 \pm 0.0%	100.0 \pm 0.0%	100.0 \pm 0.0%
	*concepts mean-step acc.	100.0 \pm 0.0%	100.0 \pm 0.0%	100.0 \pm 0.0%
	*concepts last-step acc.	100.0 \pm 0.0%	100.0 \pm 0.0%	100.0 \pm 0.0%

The concept decoder takes this as input and provides the rest of the concepts *nodesInSameSet* and *edgeInMst*, which is then passed to the output decoder to produce final $\mathbf{Y}_{kr}^{(t)}$ whether each edge is selected in the MST at step t .

The whole process is summarised in Figure 5. Yellow variables are those we supervise on. For a full list of losses used, see Appendix G.

G Algorithms and respective losses

We optimise our models based on the following losses:

- Binary cross-entropy for concepts predictions
- Binary cross-entropy for predicting termination (when applicable to the algorithm)

Additionally for each algorithm, we add:

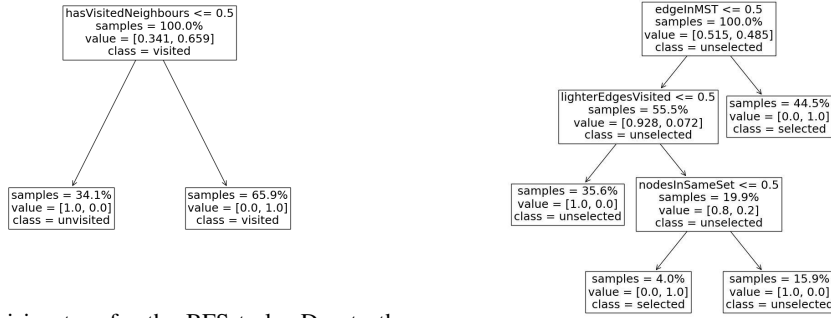
- BFS:
 - Binary cross-entropy for reachability predictions
- parallel coloring:
 - categorical cross-entropy for reachability predictions
- Kruskal’s:
 - categorical cross-entropy for predicting the position of the next selected edge
 - binary cross-entropy for predicting whether two nodes are in the same set (cf. Veličković et al. (2020a))
 - binary cross-entropy for predicting mask of the pointers to be changed (cf. Veličković et al. (2020a)) – loss from positive examples is scaled by $\frac{1}{10}$

H BFS accuracies

We present BFS accuracy in Table 6. Since BFS is relatively simple task both the baseline and the concept-bottleneck model achieve perfect accuracy.

I Node-level decision tree

The decision trees are shown in Figures 6a&7. In all cases the rules follow the logic of the original algorithm. It should be noted that the decision tree for the BFS task consists of just a single concept – why are concepts necessary then? The answer is simple – the addition of termination requires us to use the concepts *hasBeenVisited*.



(a) The decision tree for the BFS task. Due to the presence of self-loops the rule is simple – check if visited neighbours exist (same node is always present in the neighbours list).

(b) The decision tree for the Kruskal's algorithm.

Figure 6: Decision trees for BFS and Kruskal's

Table 7: BFS explanations. The model explanations match the ground-truth ones. Available concepts are *hasVisitedNeighbours* and *hasBeenVisited*. Max number of explanation occurrences is 5.

Thing to explain	Explanation	# of occurrences
n is NOT visited	$\neg hasBeenVisited(n) \wedge \neg hasVisitedNeighbours(n)$	5
n is visited	$hasVisitedNeighbours(n)$	5
continue execution	$\neg hasBeenVisited(n) \wedge hasVisitedNeighbours(n)$	5

Although not displayed in Figure 7 due to space constraints, the decision tree helped discover a 'bug' – one of the concepts (*color5Seen*) was never used, which suggests that it may not be necessary.

J Oracle model for parallel coloring

In the main body of the paper (Figure 3) we showed that quite often there were drops of accuracy for almost all concepts for the parallel coloring task. Based on more careful observation it was hypothesised that the *hasPriority* is *too critical* for the task. To prove this we plot the per-step concept accuracy of an oracle model that is always provided with the correct value for that concept. Figure 8 clearly shows the importance of the *hasPriority* concept – hardcoding its value to the ground-truth allows us to learn the task almost perfectly.

K Example explanations of algorithms

We list example explanations for BFS in Table 7 and for Kruskal's. Our observations are that for simpler algorithms with fewer possible concept combinations, such as BFS, explanations match *very closely* the ground truth – the only difference is that unvisited rule relies on $\neg hasBeenVisited$. This, however, would not break the execution – in fact the network learns to use this concept to reinforce its outputs since $\neg hasBeenVisited$ always holds when $\neg hasVisitedNeighbours$.

For the parallel coloring task, we first draw the attention to the results without using any regularisation techniques in Table 9. The main observation is that regularisation is not necessary to achieve high accuracy – formula-based metrics do not deviate largely from formulas obtained without using metrics. However, as can be seen from Table 10 when the possible combinations and outputs grow, some combinations may not be observed during training and therefore the explanations may contain unnecessary concepts (e.g., whether color 5 is seen in the explanation for uncolored node n) or sometimes not consider some concept (red in Table 10). Given the applications of the formulas still achieved very high accuracy, our hypothesis was that such algorithm require regularisation techniques, such as adding auxiliary L1 loss and pruning, in order to force the model to rely on as few concepts

Table 8: Kruskal’s explanations. The model explanations match the ground-truth ones. Available concepts are *nodesInSameSet*, *lighterEdgesVisited* (*LEV*) and *edgeInMst*. Max number of explanation occurrences is 5.

Thing to explain	Explanation	# of occurrences
e NOT in MST	$(nodesInSameSet(e) \wedge \neg edgeInMst(e)) \vee$ $(\neg LEV(e) \wedge \neg edgeInMst(e))$	5
	$(LEV(e) \wedge edgeInMst(e)) \vee$ $(LEV(e) \wedge \neg nodesInSameSet(e))$	1
e in MST	$(LEV(e) \wedge nodesInSameSet(e) \wedge edgeInMst(e)) \vee$ $(LEV(e) \wedge \neg nodesInSameSet(e) \wedge \neg edgeInMst(e))$	4

Table 9: Bottleneck model without pruning/L1 loss gives applicable accuracies – formula accuracies are as high as their counterparts. But are they as interpretable?

	Parallel coloring		
	$ V = 20$	$ V = 50$	$ V = 100$
mean-step acc.	99.55±0.17%	99.39±0.22%	99.22±0.32%
last-step acc.	99.33±0.23%	99.06±0.37%	98.95±0.52%
term. acc.	99.69±0.27%	99.41±0.21%	99.07±0.38%
formula mean-step acc.	99.55±0.18%	99.39±0.22%	99.22±0.33%
formula last-step acc.	99.31±0.18%	99.03±0.36%	98.93±0.52%
formula term. acc.	99.66±0.27%	99.41±0.21%	99.07±0.38%
concepts mean-step acc.	99.81±0.08%	99.68±0.11%	99.6±0.16%
concepts last-step acc.	99.61±0.16%	99.38±0.2%	99.27±0.27%

as possible. Results confirmed our hypothesis – Table 11 shows that 4/5 times *all* explanations are correct and only once the rule for color 2 contained an extra clause.

L Software

The code for the experiments is implemented in Python 3, relying upon open-source libraries such as Scikit-learn (BSD license) (Pedregosa et al., 2011) and Pytorch (BSD license) (Paszke et al., 2019). All the experiments have been run on an NVIDIA Titan Xp 12 GB GPU.

Table 10: Some example parallel coloring explanations for node n , *no regularisation applied*. (For brevity we omit listing every single explanation.) Class explanations do not match the ground truth and sometimes contain usage of variables, whose value can be inferred from the values of other variables. Available concepts are *isColored* (iC), *hasPriority* (hP) and *colorX Seen*, $X \in \{1, \dots, 5\}$. Max number of explanation occurrences is 5.

Thing to explain	Explanation	# of occurrences
n is not colored	$(\neg iC(n) \wedge \neg hP(n) \wedge \neg color5Seen(n)) \vee$ $(color5Seen(n) \wedge color4Seen(n) \wedge \neg iC(n) \wedge \neg hP(n)$ $\neg color1Seen(n) \wedge \neg color2Seen(n))$	1
n has color 2	$(iC(n) \wedge \neg hP(n) \wedge color1Seen(n) \wedge \neg color2Seen(n)) \vee$ $(hP(n) \wedge color1Seen(n) \wedge \neg iC(n) \wedge \neg color2Seen(n))$	2
n has color 5	$(hP(n) \wedge color2Seen(n) \wedge color3Seen(n)$ $\wedge color4Seen(n) \wedge \neg color5Seen(n) \wedge \neg iC(n)) \vee$ $(iC(n) \wedge color1Seen(n) \wedge color2Seen(n) \wedge$ $color3Seen(n) \wedge color4Seen(n) \wedge \neg color5Seen(n))$	2
n has color 5	$(hP(n) \wedge color1Seen(n) \wedge color2Seen(n) \wedge color3Seen(n)$ $\wedge color4Seen(n) \wedge \neg color5Seen(n) \wedge \neg iC(n)) \vee$ $(iC(n) \wedge color1Seen(n) \wedge color2Seen(n) \wedge$ $color3Seen(n) \wedge color4Seen(n) \wedge \neg color5Seen(n))$	3

Table 11: Parallel coloring explanations for node n , *with regularisation applied*. Class explanations are concise, consistent across seeds and very close to the ground truth. Available concepts are *isColored* (iC), *hasPriority* (hP) and *colorX Seen*, $X \in \{1, \dots, 5\}$. Max number of explanation occurrences is 5.

Thing to explain	Explanation	# of occurrences
n is not colored	$\neg iC(n) \wedge \neg hP(n)$	5
n has color 1	$(iC(n) \wedge \neg hP(n) \wedge \neg color1Seen(n)) \vee$ $(hP(n) \wedge \neg color1Seen(n) \wedge \neg iC(n))$	5
n has color 2	$(iC(n) \wedge \neg hP(n) \wedge color1Seen(n) \wedge \neg color2Seen(n)) \vee$ $(hP(n) \wedge color1Seen(n) \wedge \neg iC(n) \wedge \neg color2Seen(n)) \vee$ $(hP(n) \wedge color4Seen(n) \wedge \neg iC(n) \wedge$ $\neg color2Seen(n) \wedge \neg color3Seen(n))$	4
n has color 3	$(iC(n) \wedge \neg hP(n) \wedge color1Seen(n)$ $\wedge color2Seen(n) \wedge \neg color3Seen(n)) \vee$ $(hP(n) \wedge color1Seen(n) \wedge color2Seen(n)$ $\wedge \neg color3Seen(n) \wedge \neg iC(n))$	5
n has color 4	$(iC(n) \wedge \neg hP(n) \wedge color1Seen(n)$ $\wedge color2Seen(n) \wedge color3Seen(n) \wedge \neg color4Seen(n)) \vee$ $(hP(n) \wedge color1Seen(n) \wedge color2Seen(n)$ $\wedge color3Seen(n) \wedge \neg color4Seen(n) \wedge \neg iC(n))$	5
n has color 5	$(iC(n) \wedge \neg hP(n) \wedge color1Seen(n)$ $\wedge color2Seen(n) \wedge color3Seen(n) \wedge color4Seen(n)) \vee$ $(hP(n) \wedge color1Seen(n) \wedge color2Seen(n)$ $\wedge color3Seen(n) \wedge color4Seen(n) \wedge \neg iC(n))$	5

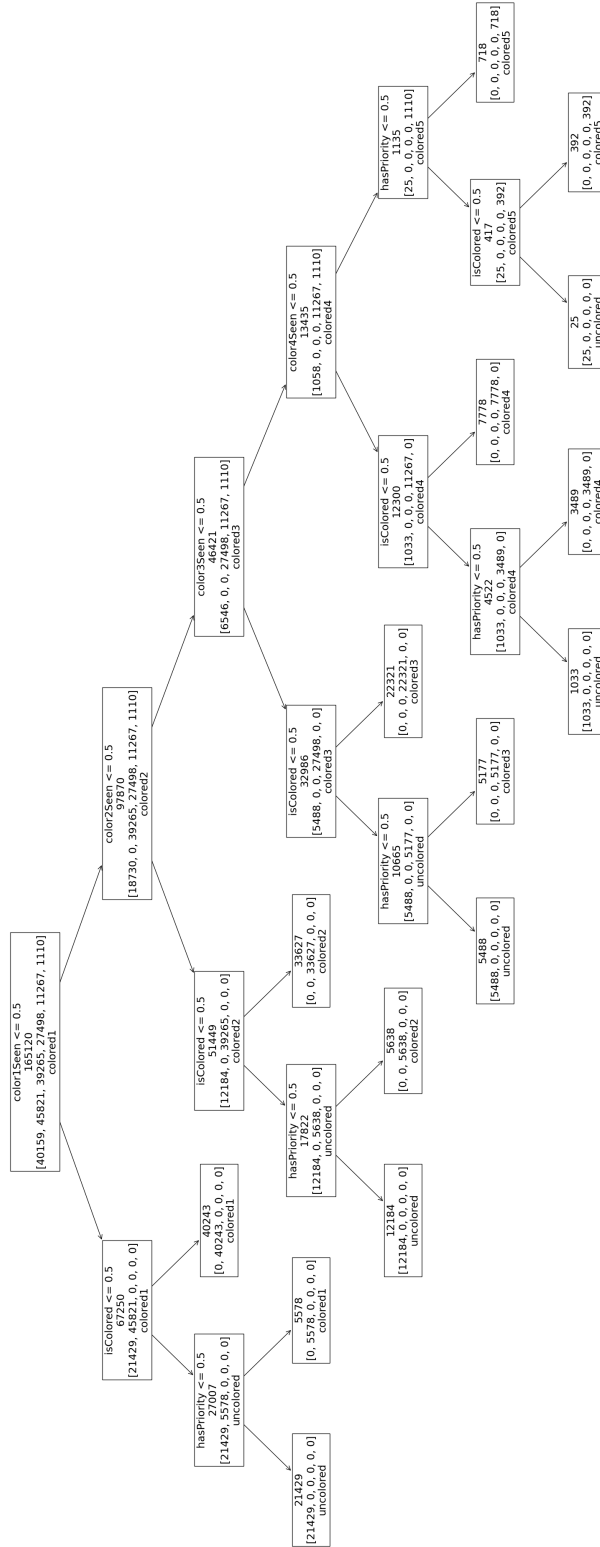


Figure 7: The decision tree for the parallel coloring task. The classifier first checks colors seen around, then priority and only then, whether a node is colored or not.

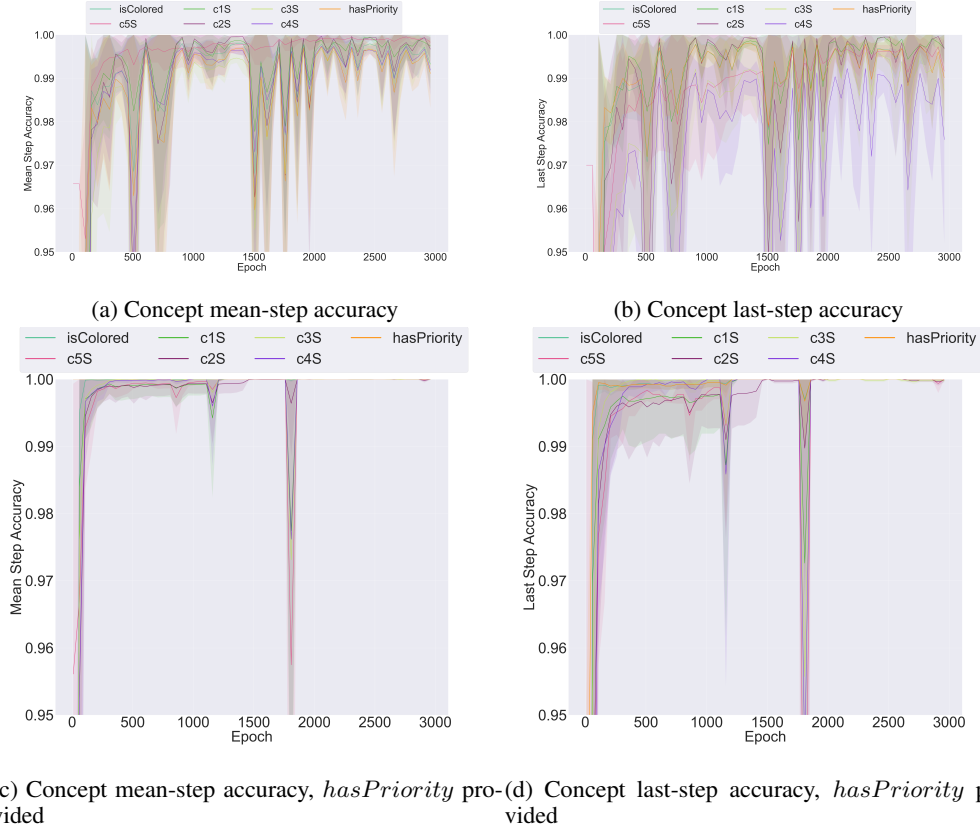


Figure 8: Concept accuracies per epoch of the parallel coloring algorithm with and without hardcoding the *hasPriority* concept. Clearly, hardcoding this concept makes our morel much more stable.