

Graph Neural Networks for Propositional Model Counting

Gaia Saveri^{1,2} and Luca Bortolussi²

¹ Department of Computer Science, University of Pisa, Italy
gaia.saveri@phd.unipi.it

² Department of Mathematics and Geoscience, University of Trieste, Italy

Abstract. Graph Neural Networks (GNNs) have been recently leveraged to solve several logical reasoning tasks. Nevertheless, counting problems such as propositional model counting (#SAT) are still mostly approached with traditional solvers. Here we tackle this gap by presenting an architecture based on the GNN framework for belief propagation (BP) of [15], extended with self-attentive GNN and trained to approximately solve the #SAT problem. We ran a thorough experimental investigation, showing that our model, trained on a small set of random Boolean formulae, is able to scale effectively to much larger problem sizes, with comparable or better performances of state of the art approximate solvers. Moreover, we show that it can be efficiently fine-tuned to provide good generalization results on different formulae distributions, such as those coming from SAT-encoded combinatorial problems.

1 Introduction

Propositional model counting (#SAT), the problem of computing the number of satisfying assignments of a given Boolean formula, is of relevant importance in computer science as it arises in many domains, such as Bayesian reasoning and combinatorial designs [7]. Nevertheless, #SAT is #P-complete, thus computationally at least as hard as NP-complete problems [25].³

For this reason, state of the art exact #SAT solvers are not capable of handling industrial-size problems, hence a number of approximate solvers (i.e. methods providing an estimation of the number of solutions) have been developed, which are able to scale to larger problem sizes.

On the other hand, in the last few years there has been an increasing interest in leveraging machine learning in general and deep learning in particular to solve logical and combinatorial reasoning tasks [18,17,5]. Graph Neural Networks (GNNs) [21] fit well in this scenario as they carry inductive biases that effectively encode combinatorial inputs, such as permutation invariance and sparsity awareness [4]. Although not always providing the same exactness guarantees of traditional solvers, the rationale of using learning-based approaches (among which GNNs) in these fields is that they can provide

³ being #P-complete the complexity class of counting problems associated with NP-complete decision problems.

a fast approximation of the computations performed by non-learned solvers, without the need of specifying distribution-specific hand-crafted heuristics.

In order to be practically exploited in combinatorial tasks, deep learning models should satisfy some desiderata [8], that we address in the experimental evaluation of this work. These properties are: scalability (i.e. the model should be able to perform well on graph instances which are from the same data generating distribution than the ones seen during training, but of larger size), generalization (i.e. the model should perform well on unseen instances), data-efficiency (i.e. the model should not require a large amount of labeled data for training, as it may be costly to obtain) and time-efficiency (i.e. the performance of the model should be close to that of optimal solvers, but faster).

In this work we investigate whether GNNs can be meaningfully applied to approximately solve the #SAT problem, an objective that, to the best of our knowledge, is only tackled in [15]. To this end, we extend the architecture presented in [15], that aims at generalizing the Belief Propagation algorithm (BP) [20] using Message Passing Neural Networks (MPNNs) [11], by augmenting the model with a self-attention mechanism.

The rationale behind keeping the algorithmic structure of BP as in [15] is that the estimate of the partition function given by BP has been proven to be a good approximation of the number of solutions of a Boolean formula [14], and also because the dynamics of MPNNs reflects the way in which information is propagated throughout the graphical model by BP. The motivation for introducing an attention mechanism inside the architecture is instead that of enhancing generalization capabilities of the model, while preserving the relational inductive biases needed in this context.

We carry out a careful and extensive experimental investigation of our proposed model, addressing in particular the already mentioned scalability and generalization properties. We show that our architecture, trained on a small set of synthetically generated random Boolean formulae, is able to scale-up to larger problem sizes, outperforming state-of-the art approximate #SAT solver. Moreover we describe a simple yet effective fine-tuning strategy, that allows the model to generalize across diverse data distributions, with only a few tens of labeled formulae required.

2 Background

2.1 Belief Propagation and #SAT

Belief Propagation [20] is an approximate inference algorithm for computing marginals of a probability distribution, exploiting the factor graph arising from its factorization. If we consider a discrete probability distribution over variables $V = \{x_1, \dots, x_n\}$ (below $x_{\mathcal{N}(f_j)}$ indicates the set of variables each factor f_j depends on):

$$P(x_1, \dots, x_n) = \frac{1}{Z} \prod_{j=1}^m f_j(x_{\mathcal{N}(f_j)}), \quad Z = \sum_{x \in V} \prod_{j=1}^m f_j(x_{\mathcal{N}(f_j)}) \quad (1)$$

then belief propagation computes an approximation of the factor marginals (also called beliefs) $\{b_j(x_{\mathcal{N}(f_j)})\}_{j=1}^m$ and of the variable marginals $\{b_i(x_i)\}_{i=1}^n$ by passing mes-

sages between neighboring nodes on the factor graph following the iterative scheme given by:

$$\begin{aligned} m_{i \rightarrow j}^{(k+1)}(x_i) &= \prod_{c \in \mathcal{N}(x_i) \setminus j} m_{c \rightarrow i}^{(k)}(x_i) \\ m_{j \rightarrow i}^{(k+1)}(x_i) &= \sum_{x_1, \dots, x_k \in \mathcal{N}(f_j) \setminus x_i} f_j(x_i, x_1, \dots, x_k) \cdot \prod_{x_v \in \mathcal{N}(f_j) \setminus x_i} m_{v \rightarrow i}^{(k)}(x_v) \end{aligned} \quad (2)$$

with initialization $m_{i \rightarrow j}^{(0)}(x_i) = 1$ for variable-to-factor messages and $m_{j \rightarrow i}^{(0)}(x_i) = f_j(x_i)$ for factor-to-variable messages. When convergence has been reached, or a pre-defined maximum number of iterations T has been performed, approximate marginals are computed. Beliefs can be used to compute a variational approximation of the partition function of the factor graph Z of Equation 1 as (below $\deg(x_i)$ denotes the number of edges incident to variable node x_i):

$$\begin{aligned} U &= - \sum_{j=1}^m \sum_{x_{\mathcal{N}(f_j)}} b_j(x_{\mathcal{N}(f_j)}) \ln f_j(x_{\mathcal{N}(f_j)}) \\ H &= - \sum_{j=1}^m \sum_{x_{\mathcal{N}(f_j)}} b_j(x_{\mathcal{N}(f_j)}) \ln b_j(x_{\mathcal{N}(f_j)}) + \sum_{i=1}^n (\deg(x_i) - 1) \sum_{x_i} b_i(x_i) \ln b_i(x_i) \\ F &= U - H \end{aligned} \quad (3)$$

where U is called Bethe average energy, H Bethe entropy and F Bethe free energy. It holds that $F \approx -\ln Z$ [6]. It is worth noting that, in order to avoid underflow errors when implementing BP, message updates of Equation 2 are usually performed in the log-space. Moreover, it has been shown that taking partial updates improves convergence of BP without changing its fixed points [14]. In detail, given a *damping* parameter $\alpha \in [0, 1]$, the following weighted combination is taken at each message passing iteration:

$$\begin{aligned} m_{i \rightarrow j}^{(k+1)}(x_i) &= \alpha \cdot m_{i \rightarrow j}^{(k+1)}(x_i) + (1 - \alpha) \cdot m_{i \rightarrow j}^{(k)} \\ m_{j \rightarrow i}^{(k+1)}(x_i) &= \alpha \cdot m_{j \rightarrow i}^{(k+1)}(x_i) + (1 - \alpha) \cdot m_{j \rightarrow i}^{(k)} \end{aligned} \quad (4)$$

It holds that propositional formulae, that without loss of generality we assume to be in Conjunctive Normal Form (CNF), can be readily translated into a factor graph. Such bipartite graph contains a factor node for every clause and a variable node for every variable in the formula, and undirected edges connecting variable nodes to the factor nodes of the clauses they appear in. If we impose that a factor node takes value 1 for variable configurations that satisfy the corresponding clause and 0 otherwise, then the partition function of this factor graph (Z of Equation 1) counts, by construction, the number of models of the input formula, i.e. the solution to the #SAT problem (implementing this

also allows to distinguish between two formulae differing only for literals’ negation, which in principle have the same factor graph representation). This intuition enables the adoption of probabilistic reasoning methods as approximate solvers for #SAT [14]. On SAT instances, BP works by iteratively passing messages between variables and clauses until a fixed point is reached. From the fixed point, an approximation of the partition function Z is computed, which serves as an approximation to the number of models of the input formula.

2.2 Graph Attention Networks

Graph Neural Networks [21] are deep learning models that address graph-related task, being natively able to process graph-structured inputs. The key idea behind GNNs is to compute a continuous representation of the nodes of the input graph that strongly depends on the structure of the graph. In particular, the Message Passing Neural Network model (MPNNs) [11], which provides a general framework for GNNs, uses a form of neural message passing, in which real-valued vector messages are exchanged between neighboring nodes, for a fixed number of iterations, in such a way that at each iteration each vertex aggregates information from its immediate neighbors.

Graph Attention Networks (GATs) [27] are a type of GNN endowed with a self-attention mechanism, that allows the network to aggregate the information coming from different nodes of the input graph putting different focus (i.e. a different weight) on some entities, and fade out the rest. More in detail, given an input graph $G = (V, E)$ and a set of continuous embeddings for the nodes of the graph $\mathbf{h} = \{\mathbf{h}_1, \dots, \mathbf{h}_{|V|}\}$, with $\mathbf{h}_i \in \mathbb{R}^d$, the graph attentional layer computes the attention coefficients between node i and node j as:

$$\alpha_{ij} = \frac{\exp(\text{LeakyReLU}(\mathbf{a}^T [\text{concat}(W\mathbf{h}_i, W\mathbf{h}_j)]))}{\sum_{k \in \mathcal{N}(i)} \exp(\text{LeakyReLU}(\mathbf{a}^T [\text{concat}(W\mathbf{h}_i, W\mathbf{h}_k)]))} \quad (5)$$

for all pairs of neighboring nodes, where $W \in \mathbb{R}^{d \times d}$ is a learnable matrix and $\mathbf{a} \in \mathbb{R}^{2d}$ is the set of parameters of a single-layer feedforward neural network. To update node representation and obtain a new set of node features $\mathbf{h}' = \{\mathbf{h}'_1, \dots, \mathbf{h}'_{|V|}\}$ the following is computed:

$$\mathbf{h}'_i = \sigma \left(\sum_{j \in \mathcal{N}(i)} \alpha_{ij} W\mathbf{h}_j \right) \quad (6)$$

being σ a non-linear differentiable function. In order to stabilize the learning process of self-attention, a multi-head attention (similar to that of [26]) is applied by replicating the operations of the layer (Equations 5 and 6) independently K times and aggregating the results as:

$$\mathbf{h}'_i = \left\| \right\|_{k=1}^K \sigma \left(\sum_{j \in \mathcal{N}(i)} \alpha_{ij}^k W^k \mathbf{h}_j \right) \quad (7)$$

where $\| \cdot \|$ denotes a feature-wise aggregation function such as concatenation, sum or average.

3 Method

The objective of this work is to tackle the #SAT problem using a GNN model as an approximate solver.

3.1 Belief Propagation Neural Networks

The starting point of our investigation is (one of the variants of) the architecture proposed in [15], that we recall here briefly. This model, called Belief Propagation Neural Network (BPNN), generalizes BP by means of GNNs, taking as input the factor graph representing a CNF SAT formula and giving as output an estimation of the logarithm of the factor graph’s partition function Z of Equation 1. Given an input factor graph $G = (V, E)$, where nodes are partitioned into factor nodes $\{f_j\}_{j=1}^m$ and variable nodes $\{x_i\}_{i=1}^n$, both factor-to-variable messages $\hat{m}_{j \rightarrow i}$ and variable-to-factor $\hat{m}_{i \rightarrow j}$ messages are initialized to 1; then the following message passing phase is performed (in the log-space) for T iterations:

$$\begin{aligned} \hat{m}_{i \rightarrow j}^{(k+1)}(x_i) &= \sum_{c \in \mathcal{N}(x_i) \setminus j} \text{MLP}_1(\hat{m}_{c \rightarrow i}^{(k)}(x_i)) \\ \hat{m}_{j \rightarrow i}^{(k+1)}(x_i) &= \text{LSE}_{x_1, \dots, x_k \in \mathcal{N}(f_j) \setminus x_i} \left(f_j(x_i, x_1, \dots, x_k) + \right. \\ &\quad \left. + \sum_{x_v \in \mathcal{N}(f_j) \setminus x_i} \text{MLP}_2(\hat{m}_{v \rightarrow i}^{(k)}(x_v)) \right) \end{aligned} \quad (8)$$

where LSE is a shorthand for the log-sum-exp function. That is, BPNN augments the standard BP message passing scheme of Equation 2 (in the log-space) by transforming the messages using MLPs.

After the message passing phase is completed, the following readout phase is executed, which outputs an estimation $\ln \hat{Z}$ of the natural logarithm of the partition function Z of the input factor graph:

$$\begin{aligned} \ln \hat{Z} &= \text{MLP}_3 \left[\text{concat}_{k=1}^T \left(\text{concat} \left(\sum_{j=1}^m b_j^{(k)}(x_{\mathcal{N}(f_j)}) \ln f_j(x_{\mathcal{N}(f_j)}), \right. \right. \right. \\ &\quad \left. \left. - \sum_{j=1}^m b_j^{(k)}(x_{\mathcal{N}(f_j)}) \ln b_j^{(k)}(x_{\mathcal{N}(f_j)}), \right. \right. \\ &\quad \left. \left. \sum_{i=1}^n (\deg(x_i) - 1) b_i^{(k)}(x_i) \ln b_i^{(k)}(x_i) \right) \right] \end{aligned} \quad (9)$$

where $b_i^{(k)}(x_i)$ and $b_j^{(k)}(x_{\mathcal{N}(f_j)})$ refer to an approximation of variable and factor beliefs at iteration k , respectively, computed as in standard belief propagation (we defer to the supplementary material, Section A for further details). Hence this final layer

takes as input a concatenation of a term representing the Bethe average energy (U of Equation 3) and two terms representing the components of the Bethe entropy (H of Equation 3) across all iterations, summed across factors within each iteration. As already mentioned in Section 2.1, damping is a standard technique for improving convergence of BP. The BPNN architecture allows for possibly applying a learned operator $\Delta : \mathbb{R}^{\sum_{i=1}^n 2 \cdot \text{deg}(x_i)} \rightarrow \mathbb{R}^{\sum_{i=1}^n 2 \cdot \text{deg}(x_i)}$ to the difference between iterations $k + 1$ and k of every factor-to-variable message, and jointly modify it, in place of the scalar multiplier α of Equation 4.

3.2 Neural Belief Propagation with Attention

Attention mechanisms, as seen in Section 2.2, are a technique that allows a deep learning model to combine input representations additively, while also enforcing permutation invariance. GNNs endowed with an attention mechanism are a promising research direction in the neuro-symbolic computing scenario [5,13,17], as they might enhance structured reasoning and efficient learning.

This is the reason for modifying the BPNN architecture by augmenting it with a GAT-style attention mechanism (i.e. analogous to that of Equations 5, 6, 7). In what follows we will refer to our architecture as BPGAT, to underline the fact that it puts together the algorithmic structure of BP and the computational formalism of GATs.

BPGAT works by taking as input a bipartite factor graph $G = (V, E)$, with factor nodes $\{f_j\}_{j=1}^m$ and variable nodes $\{x_i\}_{i=1}^n$, representing a CNF SAT formula and outputs an approximation $\ln \hat{Z}$ of the logarithm of the exact number of solutions of the input formula. As for BPNN, messages between nodes are partitioned into factor-to-variable messages and variable-to-factor messages, and computations are performed in the log-space.

Considering variable-to-factor messages, at each iteration $k + 1$ every variable node x_i contains the aggregated messages $\hat{m}_{i \rightarrow j}^{(k)}$ received from its neighborhoods at the previous iteration k , and needs to receive messages from its adjacent clause nodes, and aggregate them. Such aggregation is done using a GAT-style multi-head attention mechanism: as for the GAT attentional layer, the model computes attention coefficients α_{ij} for every factor node f_j in the neighborhood of the variable node x_i by projecting each message $\hat{m}_{j \rightarrow i}^{(k)} \in \mathbb{R}^2$ and variable node hidden representation $\hat{m}_{i \rightarrow j}^{(k)} \in \mathbb{R}^2$ via a weight matrix $W \in \mathbb{R}^{2 \times 2}$, then passing the concatenation of the projected messages through a single-layer feedforward neural network $\mathbf{a} \in \mathbb{R}^4$, then feeding the results to a LeakyReLU and finally applying softmax to normalize them across the neighborhood $\mathcal{N}(x_i)$ of x_i . More formally, attention coefficients are computed as:

$$\alpha_{ij}^{(k+1)} = \frac{\exp(\text{LeakyReLU}(\mathbf{a}^T [\text{concat}(W\hat{m}_{i \rightarrow j}^{(k)}, W\hat{m}_{j \rightarrow i}^{(k)})]))}{\sum_{k \in \mathcal{N}(x_i)} \exp(\text{LeakyReLU}(\mathbf{a}^T [\text{concat}(W\hat{m}_{i \rightarrow j}^{(k)}, W\hat{m}_{k \rightarrow i}^{(k)})]))} \quad (10)$$

The new variable-to-factor messages are then computed as a weighted sum of the incoming factor-to-variable messages, using attention coefficients of Equation 10, as:

$$\hat{m}_{i \rightarrow j}^{(k+1)}(x_i) = \sum_{c \in \mathcal{N}(x_i) \setminus j} \alpha_{ij}^{(k+1)} W \hat{m}_{c \rightarrow i}^{(k)}(x_i) \quad (11)$$

In order to stabilize the learning process, the multi-head attention mechanism of Equation 7 is used, by replicating K times the computations of Equations 10 and 11 and aggregating the results. To perform the aggregation of the outputs of each of the attention heads, concatenation has been used for the internal layers, average for the final one.

Analogous computations are performed to compute factor-to-variable messages:

$$\begin{aligned} \alpha_{ji}^{(k+1)} &= \frac{\exp(\text{LeakyReLU}(\mathbf{a}^T [\text{concat}(W \hat{m}_{j \rightarrow i}^{(k)}, W \hat{m}_{i \rightarrow j}^{(k)})]))}{\sum_{x_v \in \mathcal{N}(f_j)} \exp(\text{LeakyReLU}(\mathbf{a}^T [\text{concat}(W \hat{m}_{j \rightarrow i}^{(k)}, W \hat{m}_{i \rightarrow v}^{(k)})]))} \\ \hat{m}_{j \rightarrow i}^{(k+1)}(x_i) &= \text{LSE}_{x_1, \dots, x_k \in \mathcal{N}(f_j) \setminus x_i} \left(f_j(x_i, x_1, \dots, x_k) + \right. \\ &\quad \left. + \sum_{x_v \in \mathcal{N}(f_j) \setminus x_i} \alpha_{ji}^{(k+1)} W \hat{m}_{v \rightarrow i}^{(k)}(x_i) \right) \end{aligned} \quad (12)$$

Also for this type of messages, a multi-head attention mechanism is deployed, where the output of each head is concatenated in internal layers and averaged in the final layer.

Once this message passing phase has run for T iterations, the readout phase is executed. This consists in applying as final layer the one described in Equation 9, in order to obtain $\ln \hat{Z}$, i.e. an approximation of the natural logarithm of the number of models of the CNF SAT formula encoded by the input factor graph G .

4 Experimental Evaluation

4.1 Experimental Setting

We implemented the BPGAT architecture in Python, leveraging the PyTorch framework [19]. The model is trained to minimize the Mean Squared Error (MSE) between the natural logarithm $\ln Z$ of the true number of models of the input formula, and the output of the model $\ln \hat{Z}$.

BPGAT Training Protocol We trained the model for 1000 epochs using the Adam optimizer [12] with an initial learning rate of 10^{-4} , halving it every 200 epochs using a learning rate scheduler.

Given an input formula ϕ with n variables $\{x_i\}_{i=1}^n$ and m clauses $\{f_j\}_{j=1}^m$ and its factor graph representation $G = (V, E)$, it is preprocessed before being fed to the network in such a way that $\forall j \in \{1, \dots, m\}, \forall \{x_1, \dots, x_k\} \in \mathcal{N}(f_j), f_j(x_1, \dots, x_k) = 1$ for the assignment of $\{x_1, \dots, x_k\}$ that makes clause f_j evaluate to true, 0 otherwise, to ensure that Z of Equation 1 actually represents the count of models satisfying ϕ .

In order to compute attention coefficients of Equations 11 and 12, we used a 3-layer GAT network, having respectively 4, 4, 6 attention heads. Moreover, we used a 3-layer MLP Δ , with ReLU activation between hidden layers, to transform factor-to-variable messages, in place of a fixed-scalar damping parameter. For variable-to-factor messages a damping parameter $\alpha = 0.5$ has been used. The MLP of the final layer (Equation 9) is a 3-layer feedforward network, with ReLU non-linearity between hidden layers. The number of iterations T of the message passing scheme has been set to 5. This architectural choices have been made after performing a set of preliminary and ablation studies, whose results can be found in supplementary material, Section C.1.

BPGAT Training Data The training dataset $\mathcal{D} = \{(\phi_i, \ln Z_i)\}$ consists of a set of 1000 pairs of CNF SAT formulae ϕ_i and the logarithm $\ln Z_i$ of their true model count. Such formulae are drawn from a distribution of random formulae built in the following way:

1. Given input parameters (nv_{min}, nv_{max}) , the number of variables of the current formula ϕ is chosen as $n_{var} \sim \mathcal{U}([nv_{min}, nv_{max}])$.
2. Given input parameters (nc_{min}, nc_{max}) , the number of clauses of ϕ is chosen as $n_{cl} \sim \mathcal{U}([nc_{min}, nc_{max}])$.
3. For each clause c_j , with $j \in \{1, \dots, n_{cl}\}$, a number $k \sim 2 + \text{Bernoulli}(0.7) + \text{Geometric}(0.4)$ of variables are chosen uniformly at random from the input n_{var} variables; this leads to clause having 5 variables in average. Once the variables for the current clause c_j are chosen, each of them is negated with probability 0.5; before adding c_j to the formula ϕ , it is checked that c_j is different from all the other clauses already in the formula.
4. As soon as all the n_{cl} clauses are generated, the resulting formula ϕ is fed to Minisat [10] and checked for satisfiability: it is added to the dataset \mathcal{D} if and only if it is SAT.
5. All formulae ϕ_i in the dataset \mathcal{D} are fed to sharpSAT [24], an exact #SAT solver to generate the true count Z_i , and add its logarithm $\ln Z_i$ to the dataset.

The dataset used for training BPGAT has been generated using $(nv_{min}, nv_{max}) = (10, 30)$ and $(nc_{min}, nc_{max}) = (20, 50)$; this produced formulae having an average of 19.87 variables and 34.87 clauses.

Fine-tuning Protocol In order to allow the model to extrapolate to data distributions different from the one seen during training, without requiring a large amount of labeled data (as ground-truth labels might be costly to obtain), we use a pre-trained BPGAT as weight initializer for fine-tuning towards new formulae distributions.

In particular, we initialize the weights of our model using those of BPGAT trained for 500 epochs with the protocol described earlier in this Section and data drawn from the previously proposed random distribution, and we train it for 250 epochs with a learning rate of 10^{-6} , with only 250 labeled examples.

4.2 Results

The objective of our experiments is twofold: evaluating both BPGAT scalability and generalization capabilities. In order to assess the performance of the model, both Root Mean Squared Error (RMSE) and Mean Relative Error (MRE) metrics are reported, evaluated between the logarithm of the ground truth number of models of the input formulae $\ln Z$ and the output of the model $\ln \hat{Z}$.

As a baseline, we used ApproxMC [9,23], the state-of-the-art approximate #SAT solver, which is a randomized hashing algorithm that provides Probably Approximately Correct (PAC) guarantees.

It would have been meaningful and interesting to compare BPGAT against other guarantee-less counters, such as ApproxCount [29] or SampleApprox [28], but unfortunately we couldn't access any open-source implementation of them.

Scalability In order to assess the ability of our model to scale to larger problem sizes than the one seen during training, we generated several datasets following the procedure detailed in Section 4.1.

Table 1 shows the statistics of the datasets used in this testing phase, each containing 300 labeled instances. All datasets are much larger than the one seen during training, and in particular ‘Test 4’ contains formulae having a number of variables which is more than ten times more the one in the training set, and a number of clauses which is almost ten times more than the ones in the training set.

Table 2 shows the results obtained, in terms of RMSE and MRE, by BPGAT and ApproxMC.

It is worth noting that for all the datasets tested, BPGAT outperforms ApproxMC in terms of MRE (although not in terms of RMSE). Such higher RMSE is a consequence of few outliers with a large prediction error for BPGAT (as shown in Figure 1), while most of its predictions are close to the ground truth labels, as certified by the consistently lower MRE.

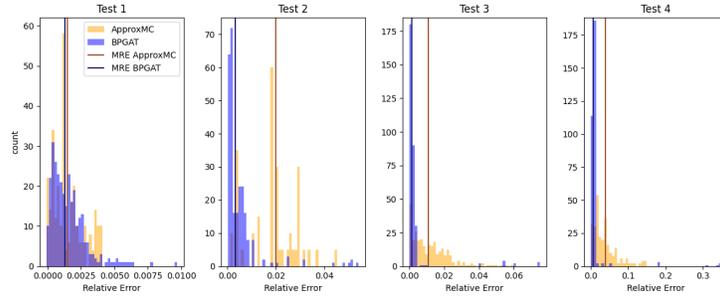
Table 1. Average number of variables, average number of clauses, average number of solutions and average time employed (in seconds) by the exact solver sharpSAT of the datasets used to test scalability.

Dataset	Avg#var	Avg#cl	Avg#sol	Avg t (s)
Test 1	61.8	76.89	1.76e+19	26.71
Test 2	60.43	143.61	6.23e+14	211.45
Test 3	124.07	75.26	1.14e+21	28.5
Test 4	377.59	275.11	7.18e+145	286.95

Out-of-distribution generalization The second set of tests we performed aims at evaluating the generalization capabilities of our model. The problem classes we perform experiments with are both SAT-encoded combinatorial problems (k -dominating set, graph

Table 2. RMSE/MRE performance of BPGAT and ApproxMC on datasets used to test scalability.

Dataset	BPGAT	ApproxMC
Test 1	0.1276/ 0.001366	0.002025 /0.001576
Test 2	0.3100/ 0.003201	0.2262 /0.02003
Test 3	0.1748/ 0.001471	0.1035 /0.01134
Test 4	1.2061/ 0.007433	0.4275 /0.04038

**Fig. 1.** Distribution of the relative errors of BPGAT and ApproxMC on datasets used to test scalability.

k -coloring, k -clique detection) and network QMR (Quick Medical Reference) problems taken from the test suite of [23].

In order to obtain SAT CNF formulae encoding the k -dominating set problem, the graph- k -coloring problem and the k -clique detection problem, we use the CNFgen tool [1]. It allows to sample graphs uniformly at random from the Erdos-Renyi random graph distribution $G(N, p)$ (specifying the two parameters $N > 0, p \in [0, 1]$) and to encode in CNF the required problem. After obtaining the CNF formulae, Minisat [10] is used to filter SAT formulae and sharpSAT [24] to obtain the number of models for each formula. Table 3 shows the statistics of the datasets used in this testing phase.

To assess the effectiveness of our fine-tuning protocol we made the following experiments (whose results are summarized in Table 4):

- We fine-tuned the model following the protocol described in Section 4.1. This is denoted as FT_BPGAT in Table 4;
- We trained, for every distribution, the model from scratch for 500 epochs, with 250 labeled examples. This is denoted as TS_BPGAT in Table 4;
- We tested BPGAT trained on random Boolean formulae (with the training protocol and the data generating procedures detailed in Section 4.1). This is denoted as BPGAT in Table 4.

As expected, the worst performance is achieved by BPGAT, as it has never been trained on formulae coming from these distributions. Interestingly, fine-tuning a model pre-trained on small random Boolean formulae (FT_BPGAT) gives better results than the same model trained on the specific dataset (TS_BPGAT). This is relevant from the

perspective of data-efficiency because small random formulae are fast to generate (differently from other distributions) and the fine-tuning phase requires only a few tens of distribution-specific labeled samples. These results are also relevant from the time-efficiency perspective, as only retraining the model for 250 epochs is needed, for each unseen data distribution.

Results of the comparison between our fine-tuned model (FT_BPGAT) and ApproxMC are shown in Table 5. It is worth observing that the performance of our architecture is comparable and in some cases outperforming that of ApproxMC (especially in terms of MRE).

Table 3. Statistics of the datasets used to test generalization. Parameters under ‘Domset’ (which stands for dominating set), ‘Color’ (which stands for graph coloring) and ‘Clique’ (which stands for clique detection) are the pair (N, p) used to generate formulae in the dataset. For all of them $k = 3$ has been used.

Dataset	Avg # var	Avg # cl	Avg # sol	Avg t (s)
Network	113.3	294.7	1.19e+20	123.9
Domset (15, 0.6)	38.43	510.15	203.3	3.76e-2
Color (10, 0.6)	65.63	294.54	572.9	3.19e-2
Clique (15, 0.5)	46.37	1145.73	102.2	3.86e-2

Table 4. RMSE/MRE performance of BPGAT when fine-tuned for the specific data distribution (FT_BPGAT), when trained on the specific data distribution (TS_BPGAT) and when trained on random formulae (BPGAT).

Dataset	FT_BPGAT	TS_BPGAT	BPGAT
Network	0.2580/0.005271	1.9334/0.04887	14.2839/0.3608
Domset	0.5508/0.04252	1.7808/0.7125	11.9190/9.5070
Color	1.2110/0.1774	1.3430/0.2046	26.1898/5.9593
Clique	0.007834/0.002475	0.01773/0.007333	1.9625/0.8983

Table 5. RMSE/MRE comparison of BPGAT fine-tuned for the specific data distributions (FT_BPGAT) and ApproxMC.

Dataset	FT_BPGAT	ApproxMC
Network	0.2580/ 0.005271	0.07619 /0.05403
Domset	0.5508/0.04252	0.08155/0.02856
Color	1.2110/0.1774	0.09426/0.04241
Clique	0.007834/0.002475	0.01113/0.04795

Data and Time-efficiency The BPGAT architecture can be claimed data-efficient, as it requires only 1000 CNF small random formulae for (pre)training. Generating such training set requires $\sim 5s$ with the procedure described in Section 4.1.

For what concerns time efficiency, the BPGAT architecture is able to process (independently on the size of the formulae), all test instances described in this work, in a maximum of 3s, without leveraging GPU acceleration. This is much less than the time needed by the exact solver sharpSAT and by the approximate solver ApproxMC, as shown in Figure 2.

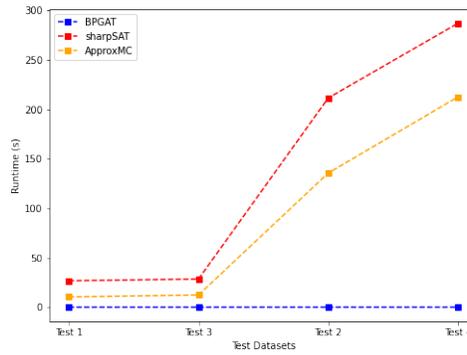


Fig. 2. Runtime comparison (in seconds) of the mean time for counting the number of satisfying assignments of a formula for the datasets used in the scalability experiments.

5 Related Work

In [2] a graph neural network model is proposed to solve the weighted disjunctive normal form counting problem (weighted #DNF). Differently from CNF #SAT, there is an approximate algorithm that admits a fully polynomial randomized approximation scheme (FPRAS), which allows the generation of hundreds of thousands of training data points, as opposed to our method which is forced to learn in a limited data regime.

A quite significant body of work has been recently developed to learn how to solve *NP*-complete combinatorial optimization problems leveraging graph neural networks [5]. Among them, methods tackling the Boolean satisfiability problem (SAT) are of particular interest for this work. They are divided into approaches leveraging GNNs as end-to-end solvers [18,22,3] (i.e. trained to output the solution directly from the input instance), and approaches in which the network is used as a computational tool to learn data-driven heuristics, in the loop of modern branch and bound solvers [30,16].

As already mentioned, our architecture builds upon the model of [15]. Besides the architectural differences that have been underlined in Section 3, the main distinction between the two models resides in the training protocol. Indeed, in [15] training data

consist in a sample of formulae drawn from the benchmarks used to test the architecture. Our model is instead trained on a set of random formulae, which are very fast to obtain, and then eventually fine-tuned towards the specific test distribution. The architectural improvements guarantee a sensibly better performance in terms of scalability and generalizability than the BPNN model, as shown in Section B of the supplementary material. We believe that the attentional layer allows the network to focus on regions of the input formulae which are more significant for the #SAT problem.

6 Conclusions

We presented BPGAT, an extension of the BPNN architecture presented in [15], which combines the algorithmic structure of belief propagation and the learning paradigm of graph attention networks. We conducted several experiments to investigate the scalability and generalization abilities of our network, showing that it is able to achieve a performance comparable to (and in some settings higher than) state-of-the-art approximate #SAT solvers, albeit the lack of any theoretical guarantees on the quality of the solution. Finally, we highlighted the efficiency of our model, both in terms of required training data and in terms of running time.

As future research directions, we will analyze the viability of extending our model to tackle weighted conjunctive normal form model counting (weighted #CNF) problems. In this scenario, a straightforward application of our model would be that of approximate probabilistic inference on Bayesian Networks, which in many cases (e.g. when solved using variational inference) comes without any statistical guarantees. It would be also interesting to generalize this approach to other counting problems, such as counting the number of independent sets of a given size in a graph, as part of a broader research program aiming at exploring the application of graph attention networks to logical reasoning tasks.

References

1. Cnfgcn: A generator of crafted benchmarks. In: Gaspers, S., Walsh, T. (eds.) Theory and Applications of Satisfiability Testing - SAT 2017 - 20th International Conference, Melbourne, VIC, Australia, August 28 - September 1, 2017, Proceedings. Lecture Notes in Computer Science, vol. 10491, pp. 464–473. Springer (2017). https://doi.org/10.1007/978-3-319-66263-3_30, https://doi.org/10.1007/978-3-319-66263-3_30
2. Abboud, R., Ceylan, İ.İ., Lukasiewicz, T.: Learning to reason: Leveraging neural networks for approximate DNF counting. In: The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020. pp. 3097–3104. AAAI Press (2020), <https://ojs.aaai.org/index.php/AAAI/article/view/5705>
3. Amizadeh, S., Matushevych, S., Weimer, M.: Learning to solve circuit-sat: An unsupervised differentiable approach. In: 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net (2019), <https://openreview.net/forum?id=BJxgz2R9t7>

4. Battaglia, P.W., Hamrick, J.B., Bapst, V., Sanchez-Gonzalez, A., Zambaldi, V.F., Malinowski, M., Tacchetti, A., Raposo, D., Santoro, A., Faulkner, R., Gülçehre, Ç., Song, H.F., Ballard, A.J., Gilmer, J., Dahl, G.E., Vaswani, A., Allen, K.R., Nash, C., Langston, V., Dyer, C., Heess, N., Wierstra, D., Kohli, P., Botvinick, M.M., Vinyals, O., Li, Y., Pascanu, R.: Relational inductive biases, deep learning, and graph networks. *CoRR* **abs/1806.01261** (2018), <http://arxiv.org/abs/1806.01261>
5. Bengio, Y., Lodi, A., Prouvost, A.: Machine learning for combinatorial optimization: A methodological tour d’horizon. *Eur. J. Oper. Res.* **290**(2), 405–421 (2021). <https://doi.org/10.1016/j.ejor.2020.07.063>, <https://doi.org/10.1016/j.ejor.2020.07.063>
6. Bethe, H.A.: Statistical Theory of Superlattices. *Proceedings of the Royal Society of London. Series A, Mathematical and Physical Sciences* **150**(871), 552–575 (1935)
7. Biere, A., Heule, M., Maaren, H.v., Walsh, T. (eds.): *Handbook of Satisfiability*, *Frontiers in Artificial Intelligence and Applications*, vol. 185. IOS Press (2009)
8. Cappart, Q., Chetelat, D., Khalil, E.B., Lodi, A., Morris, C., Velickovic, P.: Combinatorial optimization and reasoning with graph neural networks. *CoRR* **abs/2102.09544** (2021), <https://arxiv.org/abs/2102.09544>
9. Chakraborty, S., Meel, K.S., Vardi, M.Y.: Algorithmic improvements in approximate counting for probabilistic inference: From linear to logarithmic SAT calls. In: Kambhampati, S. (ed.) *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI 2016, New York, NY, USA, 9-15 July 2016*. pp. 3569–3576. IJCAI/AAAI Press (2016), <http://www.ijcai.org/Abstract/16/503>
10. Eén, N., Sörensson, N.: An extensible sat-solver. In: Giunchiglia, E., Tacchella, A. (eds.) *Theory and Applications of Satisfiability Testing, 6th International Conference, SAT 2003, Santa Margherita Ligure, Italy, May 5-8, 2003 Selected Revised Papers. Lecture Notes in Computer Science*, vol. 2919, pp. 502–518. Springer (2003). https://doi.org/10.1007/978-3-540-24605-3_37, https://doi.org/10.1007/978-3-540-24605-3_37
11. Gilmer, J., Schoenholz, S.S., Riley, P.F., Vinyals, O., Dahl, G.E.: Neural message passing for quantum chemistry. In: Precup, D., Teh, Y.W. (eds.) *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017. Proceedings of Machine Learning Research*, vol. 70, pp. 1263–1272. PMLR (2017), <http://proceedings.mlr.press/v70/gilmer17a.html>
12. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: Bengio, Y., LeCun, Y. (eds.) *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings* (2015), <http://arxiv.org/abs/1412.6980>
13. Kool, W., van Hoof, H., Welling, M.: Attention, learn to solve routing problems! In: *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net* (2019), <https://openreview.net/forum?id=ByxBFsRqYm>
14. Kroc, L., Sabharwal, A., Selman, B.: Leveraging belief propagation, backtrack search, and statistics for model counting. *Ann. Oper. Res.* **184**(1), 209–231 (2011)
15. Kuck, J., Chakraborty, S., Tang, H., Luo, R., Song, J., Sabharwal, A., Ermon, S.: Belief propagation neural networks. In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (eds.) *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual* (2020)
16. Kurin, V., Godil, S., Whiteson, S., Catanzaro, B.: Can q-learning with graph networks learn a generalizable branching heuristic for a SAT solver? In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (eds.) *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual* (2020), <https://proceedings.neurips.cc/paper/2020/hash/6d70cb65d15211726dce4c0e971e21c-Abstract.html>

17. Lamb, L.C., d'Avila Garcez, A.S., Gori, M., Prates, M.O.R., Avelar, P.H.C., Vardi, M.Y.: Graph neural networks meet neural-symbolic computing: A survey and perspective. In: Bessiere, C. (ed.) Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020. pp. 4877–4884. ijcai.org (2020). <https://doi.org/10.24963/ijcai.2020/679>, <https://doi.org/10.24963/ijcai.2020/679>
18. Li, Z., Chen, Q., Koltun, V.: Combinatorial optimization with graph convolutional networks and guided tree search. In: Bengio, S., Wallach, H.M., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (eds.) Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada. pp. 537–546 (2018), <https://proceedings.neurips.cc/paper/2018/hash/8d3bba7425e7c98c50f52ca1b52d3735-Abstract.html>
19. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E.Z., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In: Wallach, H.M., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E.B., Garnett, R. (eds.) Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada. pp. 8024–8035 (2019), <https://proceedings.neurips.cc/paper/2019/hash/bdbca288fee7f92f2bfa9f7012727740-Abstract.html>
20. Pearl, J.: Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA (1988)
21. Scarselli, F., Gori, M., Tsoi, A.C., Hagenbuchner, M., Monfardini, G.: The graph neural network model. IEEE Trans. Neural Networks **20**(1), 61–80 (2009). <https://doi.org/10.1109/TNN.2008.2005605>, <https://doi.org/10.1109/TNN.2008.2005605>
22. Selsam, D., Lamm, M., Bünz, B., Liang, P., de Moura, L., Dill, D.L.: Learning a SAT solver from single-bit supervision. In: 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net (2019), <https://openreview.net/forum?id=HJMC.iA5tm>
23. Soos, M., Meel, K.S.: BIRD: engineering an efficient CNF-XOR SAT solver and its applications to approximate model counting. In: The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019. pp. 1592–1599. AAAI Press (2019). <https://doi.org/10.1609/aaai.v33i01.33011592>, <https://doi.org/10.1609/aaai.v33i01.33011592>
24. Thurley, M.: sharpsat - counting models with advanced component caching and implicit BCP. In: Biere, A., Gomes, C.P. (eds.) Theory and Applications of Satisfiability Testing - SAT 2006, 9th International Conference, Seattle, WA, USA, August 12-15, 2006, Proceedings. Lecture Notes in Computer Science, vol. 4121, pp. 424–429. Springer (2006). https://doi.org/10.1007/11814948_38, https://doi.org/10.1007/11814948_38
25. Valiant, L.G.: The Complexity of Enumeration and Reliability Problems. SIAM J. Comput. **8**(3), 410–421 (1979), publisher: Society for Industrial and Applied Mathematics
26. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Guyon, I., von Luxburg, U., Bengio, S., Wallach, H.M., Fergus, R., Vishwanathan, S.V.N., Garnett, R. (eds.) Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA. pp. 5998–6008 (2017), <https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>

27. Velickovic, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., Bengio, Y.: Graph attention networks. In: 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net (2018), <https://openreview.net/forum?id=rJXMpikCZ>
28. Wang, J., Yin, M., Wu, J.: Two approximate algorithms for model counting. *Theor. Comput. Sci.* **657**, 28–37 (2017). <https://doi.org/10.1016/j.tcs.2016.04.047>, <https://doi.org/10.1016/j.tcs.2016.04.047>
29. Wei, W., Selman, B.: A new approach to model counting. In: Bacchus, F., Walsh, T. (eds.) *Theory and Applications of Satisfiability Testing*, 8th International Conference, SAT 2005, St. Andrews, UK, June 19-23, 2005, Proceedings. *Lecture Notes in Computer Science*, vol. 3569, pp. 324–339. Springer (2005). https://doi.org/10.1007/11499107_24, https://doi.org/10.1007/11499107_24
30. Yolcu, E., Póczos, B.: Learning local search heuristics for boolean satisfiability. In: Wallach, H.M., Larochelle, H., Beygelzimer, A., d’Alché-Buc, F., Fox, E.B., Garnett, R. (eds.) *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*. pp. 7990–8001 (2019), <https://proceedings.neurips.cc/paper/2019/hash/12e59a33dea1bf0630f46edfe13d6ea2-Abstract.html>

A Extended Background

A.1 Factor Graphs and Belief Propagation

Belief Propagation (BP) [20] is an approximate inference algorithm that leverages the factor graph arising from its input probability distribution in its computational scheme. Given a probability distribution such as the one of Equation 1, the corresponding factor graph is a bipartite graph in which nodes are partitioned into factor nodes $F = \{f_1, \dots, f_m\}$ (representing the factors of the probability distribution) and variable nodes $V = \{x_1, \dots, x_n\}$ (representing the variables of the probability distribution). There is an undirected edge between a variable node $x_i \in V$ and a factor node $f_j \in F$ if and only if f_j depends on x_i . Using this formalism, the neighbors $\mathcal{N}(f_j)$ of a factor node $f_j \in F$ are exactly all the variables appearing in f_j , and the neighbors $\mathcal{N}(x_i)$ of a variable node $x_i \in V$ are exactly all the factors that depend on that variable.

Each node passes messages to its neighbors, so that variables send messages to factor nodes, and viceversa factors send messages to variable nodes, using the iterative rules detailed in Equation 2.

It is possible to prove that BP is exact on tree factor graphs (i.e. it returns an exact solution in time linear in the number of variables n), however it can be very effective also on loopy graphs. When applied to graphs with loops, it is called Loopy Belief Propagation (LBP).

After T iterations of message passing (Equation 2), variables’ beliefs (i.e. approximate marginals) $\{b_i(x_i)\}_{i=1}^n$ are computed as:

$$\begin{aligned}
 b_i(x_i) &= \frac{1}{Z_i} \prod_{f_j \in \mathcal{N}(x_i)} m_{j \rightarrow i}^{(T)}(x_i) \\
 Z_i &= \sum_{x_i} \prod_{f_j \in \mathcal{N}(x_i)} m_{j \rightarrow i}^{(T)}(x_i)
 \end{aligned} \tag{13}$$

while factors' beliefs $\{b_j(x_{\mathcal{N}(f_j)})\}_{j=1}^m$ are computed as:

$$b_j(x_{\mathcal{N}(f_j)}) = \frac{f_j(x_{\mathcal{N}(f_j)})}{Z_j} \prod_{x_i \in \mathcal{N}(f_j)} m_{i \rightarrow j}^{(T)}(x_i) \quad (14)$$

$$Z_j = \sum_{x_{\mathcal{N}(f_j)}} f_j(x_{\mathcal{N}(f_j)}) \prod_{x_i \in \mathcal{N}(f_j)} m_{i \rightarrow j}^{(T)}$$

Using results of Equations 13 and 14, it is possible to compute the Bethe Free Energy F of Equation 3.

In order to avoid underflow errors when implementing BP, message updates of Equation 2 are usually performed in log-space, so that they read as follows:

$$\hat{m}_{i \rightarrow j}^{(k+1)}(x_i) = \sum_{c \in \mathcal{N}(x_i) \setminus j} \hat{m}_{c \rightarrow i}^{(k)}(x_i)$$

$$\hat{m}_{j \rightarrow i}^{(k+1)}(x_i) = \text{LSE}_{x_1, \dots, x_k \in \mathcal{N}(f_j) \setminus x_i} \left(f_j(x_i, x_1, \dots, x_k) + \right. \quad (15)$$

$$\left. + \sum_{x_v \in \mathcal{N}(f_j) \setminus x_i} \hat{m}_{v \rightarrow i}^{(k)}(x_v) \right)$$

where LSE stands for log-sum-exp function, defined as:

$$\text{LSE}_{x_j \setminus x_i} (f_j(x_j)) = \ln \left(\sum_{x_j \setminus x_i} \exp(f_j(x_j)) \right) \quad (16)$$

A.2 Graph Neural Networks

Graph Neural Networks (GNNs) [21] are a class of deep learning models that natively handle graph-structured data, with the goal of learning, for each vertex v in the input graph $G = (V, E)$, a state embedding $\mathbf{h}_v \in \mathbb{R}^d$, encoding information coming from neighboring nodes. The Message Passing Neural Network (MPNN) model [11] provides a general framework for GNNs, abstracting commonalities between many existing approaches in the literature. As the name suggests, the defining characteristic of MPNNs is that they use a form of neural message passing in which real-valued vector messages are exchanged between nodes (in particular between 1-hop neighboring vertices), for a fixed number of iterations.

In detail, during each message passing iteration t in a MPNN, the hidden representation $\mathbf{h}_v^{(t)}$ of each node v is updated according to the information $\mathbf{m}_v^{(t)}$ aggregated from its neighborhood $\mathcal{N}(v)$ as:

$$\mathbf{m}_v^{(t+1)} = \sum_{w \in \mathcal{N}(v)} M_t(\mathbf{h}_v^{(t)}, \mathbf{h}_w^{(t)}) \quad (17)$$

$$\mathbf{h}_v^{(t+1)} = U_t(\mathbf{h}_v^{(t)}, \mathbf{m}_v^{(t+1)})$$

where M_t is called message function and U_t is called message update function: both are arbitrary differentiable functions (typically implemented as neural networks).

After running T iterations of this message passing procedure, a readout phase is executed to compute the final embedding \mathbf{y}_v of each node v as:

$$\mathbf{y}_v = R(\{\mathbf{h}_v^{(T)} | v \in G\}) \quad (18)$$

where R (the readout function) needs again to be a differentiable function, invariant to permutations of the node states.

Hence, the intuition behind MPNNs is that at each iteration every node aggregates information from its immediate graph neighbors, so as iterations progress each node embedding contains information from further reaches of the graph.

Our model can be regarded as a MPNN taking as input the factor graph representing the input CNF SAT formula, whose message passing phase is the one detailed in Equations 11 and 12, and the readout step is that of Equation 9. It is worth noting that the distinction between factor and variable nodes makes the input graph heterogeneous. This is reflected into the MPNN architecture by using different update functions for factor and variable node embeddings, hence splitting each message passing iteration into two different phases: first every factor receives messages from its neighboring variables, then every variable receives messages from its neighboring factors.

B Comparison with BPNN

In order to experimentally show the performance improvement given by augmenting the BPNN architecture (detailed in Section 3.1) with a GAT-style attention mechanism, we also implemented and tested BPNN. The two MLPs of Equation 8 are implemented as 3-layer feedforward network, with ReLU activation between hidden layers; MLP₃ of Equation 9 and the learned operator Δ that transforms factor-to-variable messages are the same as BPGAT, detailed in Section 4.1.

Table 6 shows the results of testing both BPNN and BPGAT on the test sets described in Section 4.2, built to evaluate the scalability of the models. For all the benchmarks, BPGAT outperforms BPNN both in terms of RMSE and of MRE. A summary of these experiments is also reported on Figure 3.

Table 6. RMSE/MRE comparison between BPGAT and BPNN on datasets of Boolean random formulae.

Dataset	BPGAT	BPNN
Test 1	0.1276/0.001366	0.3140/0.004072
Test 2	0.3100/0.003201	1.3623/0.01786
Test 3	0.1748/0.001471	0.3046/0.004131
Test 4	1.2061/0.007433	2.9112/0.01681

We also performed a comparison of the generalization capabilities of BPNN and BPGAT under different training regimes, namely:

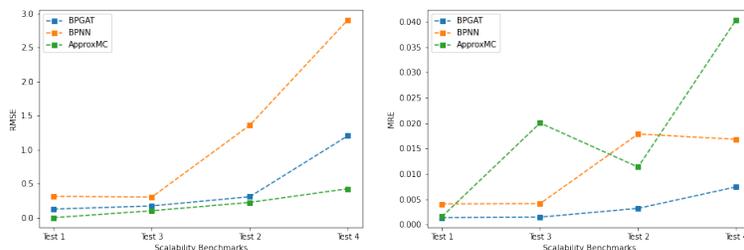


Fig. 3. Summary of the results of the experiments done on BPNN and BPGAT to test their scalability. Notation is the same of Table 6.

- We fine-tuned the two models following the protocol described in Section 4.1. Results are shown in Table 7, where FT_BPGAT and FT_BPNN denote the fine-tuned BPGAT and BPNN architectures, respectively;
- We trained, for every distribution, the two models from scratch for 500 epochs, with 250 labeled examples. Results are shown in Table 8, where TS_BPGAT and TS_BPNN denote the specifically-trained BPGAT and BPNN, respectively;
- We tested both BPGAT and BPNN trained on random Boolean formulae (with the training protocol and the data generating procedures described in Section 4.1). Results are shown in Table 9.

It is worth noting that, in all scenarios, BPGAT outperforms BPNN both in terms of RMSE and MRE. Moreover it is possible to observe that BPNN undergoes the same trend as BPGAT (Section 4.2): fine-tuning the architecture (pre-trained on random formulae) gives a better performance than training the model from scratch, using distribution-specific datasets. A summary of these experiments is also reported on Figure 4.

Table 7. RMSE/MRE comparison of BPGAT and BPNN fine-tuned for the specific distribution.

Dataset	FT_BPGAT	FT_BPNN
Network	0.2580/0.005271	0.3187/0.007469
Domset	0.5508/0.04252	1.0646/0.08392
Color	1.2110/0.1774	2.9254/0.8046
Clique	0.007834/0.002475	0.01201/0.004069

Table 8. RMSE/MRE comparison of BPGAT and BPNN trained with formulae sampled from the specific distribution.

Dataset	TS_BPGAT	TS_BPNN
Network	1.9334/0.04887	2.4358/0.0610
Domset	1.7808/0.7125	4.7885/1.6839
Color	1.3430/0.2046	4.4036/0.9088
Clique	0.01773/0.007333	0.02265/0.009737

Table 9. RMSE/MRE comparison of BPGAT and BPNN trained on random formulae.

Dataset	BPGAT	BPNN
Network	14.2839/0.3608	33.0484/0.8328
Domset	11.9190/9.5070	38.7884/13.3770
Color	26.1898/5.9593	36.1179/7.8826
Clique	1.9625/0.8983	5.2792/2.4177

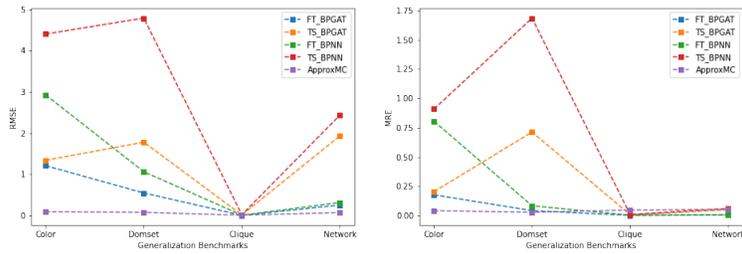


Fig. 4. Summary of the results of the experiments done on BPNN and BPGAT to test their generalization capabilities. Notation is the same of Tables 8 and 9.

C Configuration and Ablation Studies

C.1 BPGAT Parameter Tuning

Before setting BPGAT architecture’s parameters to the ones detailed in Section 4.1, we performed several tuning experiments.

Learning the Damping Parameter In order to assess whether it is better to use a fixed scalar parameter α or a learned operator Δ as damping parameter of Equation 4, we tried all possible configurations (each configuration is trained as specified in Section 4.1).

In Table 10, BPGAT refers to the architecture with damping parameter computed by an MLP Δ for factor-to-variable messages, and fixed to 0.5 for variable-to-factor messages; BPGAT_VF refers to the architecture with damping parameter computed by an MLP Δ for variable-to-factor messages, and fixed to 0.5 for factor-to-variable messages; BPGAT_ALL refers to the model trained with an MLP Δ for transforming both kind of messages; BPGAT_NONE refers to the model trained with both damping parameter fixed to 0.5. Results justify our choice of using a learned operator Δ only for updating factor-to-variable messages.

Table 10. RMSE/MRE comparison of BPGAT trained with different configurations of the damping parameter.

Dataset	BPGAT	BPGAT_VF	BPGAT_ALL	BPGAT_NONE
Test 1	0.1276/0.001366	0.2351/0.002762	0.3157/0.003930	0.4158/0.006443
Test 2	0.3100/0.003201	0.8525/0.01001	0.9493/0.01601	0.9238/0.01555
Test 3	0.1748/0.001471	0.2109/0.002998	0.2589/0.003041	0.3013/0.003344

Number of Message Passing Iterations Another parameter that we explored is the number T of message passing rounds, i.e. the number of times Equations 11 and 12 are iterated before the readout phase (Equation 9) is executed. Results are shown in Table 11: the model achieves the highest performance when $T = 5$.

Table 11. RMSE/MRE comparison for different number of message passing iterations T . BPGAT, BPGAT_10 and BPGAT_15 refer to the model trained with $T = 5, 10, 15$, respectively.

Dataset	BPGAT	BPGAT_10	BPGAT_15
Test 1	0.1276/0.001366	0.1684/0.003415	0.1775/0.003694
Test 2	0.3100/0.003201	0.3751/0.005622	0.5328/0.006294
Test 3	0.1748/0.001471	0.2555/0.003165	0.2784/0.005979

C.2 Ablation Studies on GAT Layers

To assess the relevance of using a GAT-style attention mechanism in transforming both factor-to-variable and variable-to-factor messages, we performed the following ablation studies:

- Hybrid BP-BPGAT, results are shown in Table 12: we tested the model transforming factor-to-variable messages using GAT (Equation 12) and variable-to-factor messages using BP (Equation 2), this is denoted as FVGAT-VFNONE in Table 12; conversely, we tested the model transforming variable-to-factor messages using BP (Equation 2) and factor-to-variable messages using GAT (Equation 11), this is denoted as FVNONE-VFGAT in Table 12;
- Hybrid BPNN-BPGAT, results are shown in Table 13: we tested the model transforming factor-to-variable messages using GAT (Equation 12) and variable-to-factor messages using BPNN’s corresponding update (Equation 8), this is denoted as FVGAT-VFMLP in Table 13; conversely, we tested the model transforming variable-to-factor messages using BPNN’s update (Equation 8) and factor-to-variable messages using GAT (Equation 11), this is denoted as FVMLP-VFGAT in Table 13;

Results highlight the fact that transforming both factor-to-variable and variable-to-factor messages using a GAT attention mechanism is beneficial for the overall performance of the model.

Table 12. RMSE/MRE results of the ablation studies on the attention mechanism.

Dataset	BPGAT	FVGAT-VFNONE	FVNONE-VFGAT
Test 1	0.1276/0.001366	0.1568/0.0015721	0.1786/0.001740
Test 2	0.3100/0.003201	0.4292/0.004640	0.5870/0.006392
Test 3	0.1748/0.001471	0.2009/0.001964	0.2699/0.002197

Table 13. RMSE/MRE results of the ablation studies between message transformation performed by BPGAT and BPNN.

Dataset	BPGAT	FVGAT-VFMLP	FVMLP-VFGAT
Test 1	0.1276/0.001366	0.3613/0.004713	0.4056/0.005256
Test 2	0.3100/0.003201	1.5787/0.02002	1.7065/0.02211
Test 3	0.1748/0.001471	0.2751/0.004689	0.2902/0.004872