# Hybrid Models for Learning to Branch

**Prateek Gupta**[*]
University of Oxford
The Alan Turing Institute
pgupta@robots.ox.ac.uk

**Maxime Gasse**
Mila, Polytechnique Montréal
maxime.gasse@polymtl.ca

**Elias B. Khalil**
University of Toronto
khalil@mie.utoronto.ca

**M. Pawan Kumar**
University of Oxford
pawan@robots.ox.ac.uk

**Andrea Lodi**
CERC, Polytechnique Montréal
andrea.lodi@polymtl.ca

**Yoshua Bengio**
Mila, Université de Montréal
yoshua.bengio@mila.quebec

## Abstract

A recent Graph Neural Network (GNN) approach for learning to branch has been shown to successfully reduce the running time of branch-and-bound (B&B) algorithms for Mixed Integer Linear Programming (MILP). While the GNN relies on a GPU for inference, MILP solvers are purely CPU-based. This severely limits its application as many practitioners may not have access to high-end GPUs. In this work, we ask two key questions. First, in a more realistic setting where only a CPU is available, is the GNN model still competitive? Second, can we devise an alternate computationally inexpensive model that retains the predictive power of the GNN architecture? We answer the first question in the negative, and address the second question by proposing a new hybrid architecture for efficient branching on CPU machines. The proposed architecture combines the expressive power of GNNs with computationally inexpensive multi-layer perceptrons (MLP) for branching. We evaluate our methods on four classes of MILP problems, and show that they lead to up to 26% reduction in solver running time compared to state-of-the-art methods without a GPU, while extrapolating to harder problems than it was trained on. The code for this project is publicly available at https://github.com/pg2455/Hybrid-learn2branch.

## 1 Introduction

Mixed-Integer Linear Programs (MILPs) arise naturally in many decision-making problems such as auction design [1], warehouse planning [13], capital budgeting [14] or scheduling [15]. Apart from a linear objective function and linear constraints, some decision variables of a MILP are required to take integral values, which makes the problem NP-hard [35].

Modern mathematical solvers typically employ the B&B algorithm [29] to solve general MILPs to global optimality. While the worst-case time complexity of B&B is exponential in the size of the problem [38], it has proven efficient in practice, leading to wide adoption in various industries. At a high level, B&B adopts a divide-and-conquer approach that consists in recursively partitioning the original problem into a tree of smaller sub-problems, and solving linear relaxations of the sub-problems until an integral solution is found and proven optimal.

---

[*]The work was done during an internship at Mila and CERC. Correspondence to: <pgupta@robots.ox.ac.uk>

Despite its apparent simplicity, there are many practical aspects that must be considered for B&B to perform well [2]; such decisions will affect the search tree, and ultimately the overall running time. These include several decision problems [33] that arise during the execution of the algorithm, such as *node selection*: which sub-problem do we analyze next?; and *variable selection* (a.k.a. branching): which decision variable must be used (branched on) to partition the current sub-problem? While such decisions are typically made using hard-coded expert heuristics which are implemented in modern solvers, more and more attention is given to statistical learning approaches for replacing and improving upon those heuristics [23; 5; 26; 17; 39]. An extensive review of different approaches at the intersection of statistical learning and combinatorial optimization is given in Bengio et al. [7].

Recently, Gasse et al. [17] proposed to tackle the variable selection problem in B&B using a Graph Neural Network (GNN) model. The GNN exploits the bipartite graph formulation of MILPs together with a shared parametric representation, thus allowing it to model problems of arbitrary size. Using imitation learning, the model is trained to approximate a very good but computationally expensive "expert" heuristic named *strong branching* [6]. The resulting branching strategy is shown to improve upon previously proposed approaches for branching on several MILP problem benchmarks, and is competitive with state-of-the-art B&B solvers. We note that one limitation of this approach, with respect to general B&B heuristics, is that the resulting strategy is only tailored to the class of MILP problems it is trained on. This is very reasonable in our view, as practitioners usually only care about solving very specific problem types at any time.
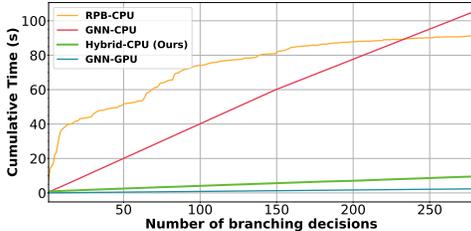


Figure 1: Cumulative time cost of different branching policies: (i) the default internal rule RPB of the SCIP solver; (ii) a GNN model (using a GPU or a CPU); and (iii) our hybrid model. Clearly the GNN model requires a GPU for being competitive, while our hybrid model does not. (Measured on a capacitated facility location problem, medium size).

While the GNN model seems particularly suited for learning branching strategies, one drawback is a high computational cost for inference, i.e., choosing the branching variable at each node of the B&B tree. In Gasse et al. [17], the authors use a high-end GPU card to speed-up the GNN inference time, which is a common practice in deep learning but is somewhat unrealistic for MILP practitioners. Indeed, commercial MILP solvers rely solely on CPUs for computation, and the GNN model from Gasse et al. [17] is not competitive on CPU-only machines, as illustrated in Figure 1. There is indeed a trade-off between the quality of the branching decisions made and the time spent obtaining those decisions. This trade-off is well-known in MILP community [2], and has given rise to carefully balanced strategies designed by MILP experts, such as *hybrid branching* [3], which derive from the computationally expensive *strong branching* heuristic [6].

In this paper, we study the time-accuracy trade-off in learning to branch with the aim of devising a model that is both computationally inexpensive and accurate for branching. To this end, we propose a hybrid architecture that uses a GNN model only at the root node of the B&B tree and a weak but fast predictor, such as a simple Multi-Layer Perceptron (MLP), at the remaining nodes. In doing so, the weak model is enhanced by high-level structural information extracted at the root node by the GNN model. In addition to this new hybrid architecture, we experiment and evaluate the impact of several variants to our training protocol for learning to branch, including: (i) end-to-end training [19; 8], (ii) knowledge distillation [25], (iii) auxiliary tasks [30], and (iv) depth-dependent weighting of the training loss for learning to branch, an idea originally proposed by He et al. [23] in the context of node selection.

We evaluate our approach on large-scale MILP instances from four problem families: Capacitated Facility Location, Combinatorial Auctions, Set Covering, and Independent Set. We demonstrate empirically that our combination of hybrid architecture and training protocol results in state-of-the-art performance in the realistic setting of a CPU-restricted machine. While we observe a slight decrease in the predictive performance of our model with respect to the original GNN from [17], its reduced computational cost still allows for a reduction of up to 26% in overall solving time on all the evaluated benchmarks compared to the default branching strategy of the modern open-source solver SCIP [20]. Also, our hybrid model preserves the ability to extrapolate to harder problems than trained on, as the original GNN model.

## 2 Related Work

Finding a branching strategy that results in the smallest B&B tree for a MILP is at least as hard–and possibly much harder–than solving the MILP with any strategy. Still, small trees can be obtained by using a computationally expensive heuristic named *strong branching* (SB) [6; 31]. The majority of the research efforts in variable selection are thus aimed at matching the performance of SB through faster approximations, via cleverly handcrafted heuristics such as *reliability pseudocost branching* [4], and recently via machine learning[2] [5; 26; 17]. We refer to [33] for an extensive survey of the topic.

Alvarez et al. [5] and Khalil et al. [26] showed that a fast discriminative classifier such as extremely randomized trees [18] or support vector machines [24] on hand-designed features can be used to mimic SB decisions. Subsequently, Gasse et al. [17] and Zarpellon et al. [39] showed the importance of representation learning for branching. Our approach, in some sense, combines the superior representation framework of Gasse et al. [17] with the computationally cheaper framework of Khalil et al. [26]. Such hybrid architectures have been successfully used in ML problems such as visual reasoning [36], style-transfer [11], natural language processing [37; 10], and speech recognition [27].

## 3 Preliminaries

Throughout this paper, we use boldface for vectors and matrices. A MILP is a mathematical optimization problem that combines a linear objective function, a set of linear constraints, and a mix of continuous and integral decision variables. It can be written as:

$$\arg \min_{\mathbf{x}} \mathbf{c}^\mathsf{T}\mathbf{x}, \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} \leq \mathbf{b}, \quad \mathbf{x} \in \mathbb{Z}^p \times \mathbb{R}^{n-p},$$

where $\mathbf{c} \in \mathbb{R}^n$ denotes the cost vector, $\mathbf{A} \in \mathbb{R}^{m \times n}$ the matrix of constraint coefficients, $\mathbf{b} \in \mathbb{R}^m$ is the vector of constant terms of the constraints, and there are $p$ integer variables, $1 \leq p \leq n$ .

The B&B algorithm can be described as follows. One first solves the linear program (LP) relaxation of the MILP, obtained by disregarding the integrality constraints on the decision variables. If the LP solution $\mathbf{x}^\star$ satisfies the MILP integrality constraints, or is worse than a known integral solution, then there is no need to proceed further. If not, then one divides the MILP into two sub-MILPs. This is typically done by picking an integral decision variable that has a fractional value, $i \in \mathcal{C} = \{i \mid x_i^\star \notin \mathbb{Z}, i \leq p\}$, and create two sub-MILPs with additional constraints $x_i \leq \lfloor x_i^\star \rfloor$ and $x_i \geq \lceil x_i^\star \rceil$, respectively. The decision variable $i$ that is used to partition the feasible region is called the *branching variable*, while $\mathcal{C}$ denotes the *branching candidates*. The second step is to select one of the leaves of the tree, and repeat the above steps until all leaves have been processed[3].

In this work, we refer to the first node processed by B&B as the *root node*, which contains the original MILP, and all subsequent nodes containing a local MILP as *tree nodes*, whenever the distinction is required. Otherwise we refer to them simply as nodes.

## 4 Methodology

As mentioned earlier, computationally heavy GNNs can be prohibitively slow when used for branching on CPU-only machines. In this section we describe our hybrid alternative, which combines the superior inductive bias of a GNN at the root node with a computationally inexpensive model at the tree nodes. We also discuss various enhancements to the training protocol, in order to enhance the performance of the learned models.

### 4.1 Hybrid architecture

A variable selection strategy in B&B can be seen as a scoring function $f$ that outputs a score $s_i \in \mathbb{R}$ for every branching candidate. As such, $f$ can be modeled as a parametric function, learned by ML. Branching then simply involves selecting the highest-scoring candidate according to $f$:

$$i_f^\star = \arg \max_{i \in \mathcal{C}} \mathbf{s}_i$$

---

[2]Interestingly, early works on ML methods for branching can be traced back to 2000 [2, Acknowledgements].
[3]For a more involved description of B&B, the reader is referred to Achterberg [2].

Table 1: Various functional forms $f$ considered for variable selection ($\odot$ denotes Hadamard product).

| | Data extraction | Computational Cost | Decision Function | |
|---|---|---|---|---|
| GNN [17] | expensive | Expensive | $\mathbf{s}$ = | $\mathrm{GNN}(\mathbf{G})$ |
| MLP | cheap | Moderate | $\mathbf{s}$ = | $\mathrm{MLP}(\mathbf{X})$ |
| CONCAT | hybrid | Moderate | $\boldsymbol{\Psi}$ = $\mathbf{s}$ = | $\mathrm{GNN}(\mathbf{G}^0)$ $\mathrm{MLP}([\boldsymbol{\Psi}, \mathbf{X}])$ |
| FiLM [36] | hybrid | Moderate | $\boldsymbol{\gamma}, \boldsymbol{\beta}$ = $\mathbf{s}$ = | $\mathrm{GNN}(\mathbf{G}^0)$ $\mathrm{FiLM}(\boldsymbol{\gamma}, \boldsymbol{\beta}, \mathrm{MLP}(\mathbf{X}))$ |
| HyperSVM | hybrid | Cheapest | $\mathbf{W}$ = $\mathbf{s}$ = | $\mathrm{GNN}(\mathbf{G}^0)$ $(\mathbf{W} \odot \mathbf{X})\mathbf{1}$ |
| HyperSVM-FiLM | hybrid | Cheapest | $\boldsymbol{\gamma}, \boldsymbol{\beta_1}, \boldsymbol{\beta_2}$ = $\mathbf{s}$ = | $\mathrm{GNN}(\mathbf{G}^0)$ $\boldsymbol{\beta_2}^\mathsf{T} \max(0, \boldsymbol{\beta_1} \odot \mathbf{X} + \boldsymbol{\gamma})$ |

We consider two forms of node representations for machine learning models: (i) a graph representation $\mathbf{G} \in \mathcal{G}$, such as the variable-constraint bipartite graph of Gasse et al. [17], where $\mathbf{G} = (\mathbf{V}, \mathbf{E}, \mathbf{C})$, with $\mathbf{V} \in \mathbb{R}^{n \times d_1}$ variable features, $\mathbf{E} \in \mathbb{R}^{n \times m \times d_2}$ edge features, and $\mathbf{C} \in \mathbb{R}^{m \times d_3}$ constraint features; and (ii) branching candidate features $\mathbf{X} \in \mathbb{R}^{|\mathcal{C}| \times d_4}$, such as those from Khalil et al. [26], which is cheaper to extract than the first representation. For convenience, we denote by $\mathcal{X}$ the generic space of the branching candidate features, and by $\mathbf{G}^0$ the graph representation of the root node. The various features $d_i$ used in this work are detailed in the supplementary materials.

In B&B, structural information in the tree nodes, $\mathbf{G}$, shares a lot of similarity with that of the root node, $\mathbf{G}^0$. Extracting, but also processing that information at every node is an expensive task, which we will try to circumvent. The main idea of our hybrid approach is then to succinctly extract the relevant structural information only once, at the root node, with a parametric model $\mathrm{GNN}(\mathbf{G}^0; \boldsymbol{\theta})$. We then combine in the tree nodes this preprocessed structural information with the cheap candidate features $\mathbf{X}$, using a hybrid model $f := \mathrm{MLP}_{\mathrm{HYBRID}}(\mathrm{GNN}(\mathbf{G}^0; \boldsymbol{\theta}), \mathbf{X}; \boldsymbol{\phi})$. By doing so, we hope that the resulting model will approach the performance of an expensive but powerful $f := \mathrm{GNN}(\mathbf{G})$, at almost the same cost as an inexpensive but less powerful $f := \mathrm{MLP}(\mathbf{X})$. Figure 2 illustrates the differences between those approaches, in terms of data extraction.

For an exhaustive coverage of the computational spectrum of hybrid models, we consider four ways to enrich the feature space of an MLP via a GNN's output, sumarrized in Table 1. In CONCAT, we concatenate the candidate's *root representations* $\boldsymbol{\Psi}$ with the features $\mathbf{X}$ at a node. In FiLM [36], we generate film parameters $\boldsymbol{\gamma}$, $\boldsymbol{\beta}$ from the GNN, for each candidate, which are further used to modulate the hidden layers of the MLP. In details, if $\boldsymbol{h}$ is the intermediate representation of the MLP, it gets linearly modulated as $\boldsymbol{h} \leftarrow \boldsymbol{\beta} \cdot \boldsymbol{h} + \boldsymbol{\gamma}$. While both the above architectures have similar computational complexity, it has been shown that FiLM subsumes the CONCAT architecture [12]. On the other the end of the spectrum lie the most inexpen-
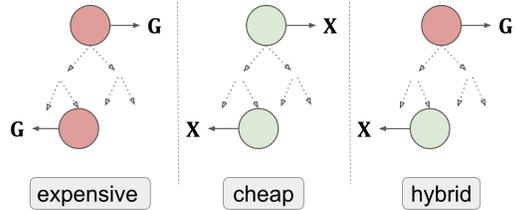


Figure 2: Data extraction strategies: bipartite graph representation $\mathbf{G}$ at every node (expensive); candidate variable features $\mathbf{X}$ at every node (cheap); bipartite graph at the root node and variable features at tree node (hybrid).

sive hybrid architectures, HyperSVM and HyperSVM-FiLM. HyperSVM is inspired by HyperNetworks [22], and simply consists in a multi-class Support Vector Machine (SVM), whose parameters are predicted by the root GNN. We chose a simple linear discriminator, for a minimal computational cost. Finally, in HyperSVM-FiLM we increase the expressivity of HyperSVM with the help of modulations, similar to that of FiLM.

## 4.2 Training Protocol

We use *strong branching* decisions as ground-truth labels for imitation learning, and collect observations of the form $(\mathbf{G}^0, \mathbf{G}, \mathbf{X}, i_{SB}^\star)$. Thus, the data used for training the model is $\mathcal{D} =$

$\{(\mathbf{G}_k^0, \mathbf{G}_k, \mathbf{X}_k, i_{SB,k}^\star), k = 1, 2, ..., N\}$. We treat the problem of identifying $i_{SB,k}^\star$ as a classification problem, such that $i^* = \arg\max_{i \in \mathcal{C}} f(\mathbf{G}^0, \mathbf{X})$ is the target outcome. Considering $\mathcal{D}$ as our ground-truth, our objective (1) is to minimize the cross-entropy loss $l$ between $f(\mathbf{G}^0, \mathbf{X}) \in \mathcal{R}^{|\mathcal{C}|}$ and a one-hot vector with one at the target:

$$\mathcal{L}(\mathcal{D}; \boldsymbol{\theta}, \boldsymbol{\phi}) = \frac{1}{N} \sum_{k=1}^{N} l(f(\mathbf{G}_k^0, \mathbf{X}_k; \boldsymbol{\theta}, \boldsymbol{\phi}), i_{SB,k}^\star). \tag{1}$$

Performance on unseen instances, or generalization, is of utmost importance when the trained models are used on bigger instances. The choices in training the aforementioned architectures influence this ability. In this section, we discuss four such important choices that lead to better generalization.

### 4.2.1 End-to-end Training (e2e)

A $\text{GNN}_{\hat{\boldsymbol{\theta}}}$, with pre-trained parameters $\hat{\boldsymbol{\theta}}$, is obtained using the procedure described in Gasse et al. [17]. We use this pre-trained GNN to extract variable representations at the root node and use it as an input to the MLPs at a tree node (pre). This results in a considerable performance boost over plain "salt-of-the-earth" MLP used at all tree nodes. However, going a step further, an end-to-end (e2e) approach involves training the GNN and the MLP together by backpropagating the gradients from a tree node to the root node. In doing so, e2e training aligns the variable representations at the root node with the prediction task at a tree node. At the same time, it is not obvious that it should result in a stable learning behavior because the parameters for GNN need to adapt to various tree nodes. Our experiments explore both pre-training (pre) and end-to-end training (e2e), namely:

$$\text{(pre)} \quad \boldsymbol{\phi}^* = \arg\min_{\boldsymbol{\phi}} \mathcal{L}(\mathcal{D}; \hat{\boldsymbol{\theta}}, \boldsymbol{\phi}), \quad \text{(e2e)} \quad \boldsymbol{\phi}^*, \boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}, \boldsymbol{\phi}} \mathcal{L}(\mathcal{D}; \boldsymbol{\theta}, \boldsymbol{\phi}). \tag{2}$$

### 4.2.2 Knowledge Distillation (KD)

Using the outputs of a pre-trained expert model as a soft-target for training a smaller model has been successfully used in model compression [25]. In this way, one aims to learn an inexpensive model that has the same generalization power as an expert. Thus, for better generalization, instead of training our hybrid architectures with cross-entropy [21] on ground-truth hard-labels, we study the effect of training with KL Divergence [28] between the outputs of a pre-trained GNN and a hybrid model, namely:

$$\text{(KD)} \quad \mathcal{L}_{KD}(\mathcal{D}; \boldsymbol{\theta}, \boldsymbol{\phi}) = \frac{1}{N} \sum_{k=1}^{N} \text{KL}(f(\mathbf{G}_k^0, \mathbf{X}_k; \boldsymbol{\theta}, \boldsymbol{\phi}), \text{GNN}_{\hat{\boldsymbol{\theta}}}(\mathbf{G_k})). \tag{3}$$

### 4.2.3 Auxiliary Tasks (AT)

An inductive bias, such as GNN, encodes a prior on the way to process raw input. Auxiliary tasks, on the other hand, inject priors in the model through additional learning objectives, which are not directly linked to the main task. These tasks are neither related to the final output nor do they require additional training data. One such auxiliary task is to maximize the diversity in variable representations. The intuition is that very similar representations lead to very close MLP score predictions, which is not useful for branching.

We minimize a pairwise loss function that ensures maximum separation between the variable representations projected on a unit hypersphere. We consider two types of objectives for this: (i) Euclidean Distance (ED), and (ii) Minimum Hyperspherical Energy (MHE) [32], inspired from the well-known Thomson problem [16] in Physics. While ED separates the representations in the Euclidean space on the hypersphere, MHE ensures uniform distribution over the hypersphere. Denoting $\hat{\boldsymbol{\psi}}_i$ as the variable representation for the variable $i$ projected on a unit hypersphere and $e_{ij} = ||\hat{\boldsymbol{\psi}}_i - \hat{\boldsymbol{\psi}}_j||_2$ as the Euclidean distance between the representations for the variables $i$ and $j$, our new objective function is given as $\mathcal{L}_{AT}(\mathcal{D}; \boldsymbol{\theta}, \boldsymbol{\phi}) = \mathcal{L}(\mathcal{D}; \boldsymbol{\theta}, \boldsymbol{\phi}) + g(\boldsymbol{\Psi}; \boldsymbol{\theta})$, where

$$\text{(ED)} \quad g(\boldsymbol{\Psi}; \boldsymbol{\theta}) = \frac{1}{N^2} \sum_{i,j=1}^{N} e_{ij}^2, \quad \text{(MHE)} \quad g(\boldsymbol{\Psi}; \boldsymbol{\theta}) = \frac{1}{N^2} \sum_{i,j=1}^{N} \frac{1}{e_{ij}}. \tag{4}$$

#### 4.2.4 Loss Weighting Scheme

The problem of distribution shift is unavoidable in a sequential process like B&B. A suboptimal branching decision at a node closer to the root node can have worse impact on the size of the B&B tree as compared to when such a decision is made farther from it. In such situations, one can use depth (possibly normalized) as a feature, but the generalization on bigger instances is a bit unpredictable as the distribution of this feature might be very different from that observed in the training set. Thus, we experimented with different depth-dependent formulations for weighting the loss at any node. Denoting $z_i$ as the depth of a tree node $i$ relative to the depth of the tree, we weight the loss at different tree nodes by $w(z_i)$, making our objective function as

$$\mathcal{L}(\mathcal{D}; \boldsymbol{\theta}, \boldsymbol{\phi}) = \frac{1}{N} \sum_{k=1}^{N} w(z_k) \cdot l(f(\mathbf{G}_k^0, \mathbf{X}_k; \boldsymbol{\theta}, \boldsymbol{\phi}), i_{SB,k}^{\star}). \tag{5}$$

Specifically, we considered 5 different weighting functions such that all of them have the same end points, i.e., $w(0) = 1.0$ at the root node and $w(1) = e^{-0.5}$ at the deepest node. Different functions were chosen depending on their intermediate behaviour in between these two points. We experimented with exponential, linear, quadratic and sigmoidal decay behavior of these functions. Table 3 lists various functions and their mathematical forms considered in our experiments.

## 5 Experiments

We follow the experimental setup of Gasse et al. [17], and evaluate each branching strategy across four different problem classes, namely Capacitated Facility Location, Minimum Set Covering, Combinatorial Auctions, and Maximum Independent Set. Randomly generated instances are solved offline using SCIP [20] to collect training samples of the form $(\mathbf{G}^0, \mathbf{G}, \mathbf{X}, i_{SB}^{\star})$. We leave the description of the data collection and training details of each model to the supplementary materials.

**Evaluation.** As in Gasse et al. [17], our evaluation instances are labeled as small, medium, and big based on the size of underlying MILP. Small instances have the same size as those used to generate the training datasets, and thus match the training distribution, while instances of increasing size allows us to measure the generalization ability of the trained models. Each scenario uses 20 instances, solved using 3 different seeds to account for solver variability. We report standard metrics used in the MILP community for benchmarking B&B solvers: (i) Time: 1-shifted geometric mean[4] of running times in seconds, including the running times for unsolved instances, (ii) Nodes: hardware-independent 1-shifted geometric mean of B&B node count of the instances solved by all branching strategies , and (iii) Wins: number of times each branching strategy resulted in the fastest solving time, over total number of solved instances. All branching strategies are evaluated using the open-source solver SCIP [20] with a time limit of 45 minutes, and cutting planes are allowed only at the root node.

**Baselines.** We compare our hybrid model to several non-ML baselines, including SCIP's default branching heuristic *Reliability Pseudocost Branching* (RPB), the "gold standard" heuristic *Full Strong Branching* (FSB), and a very fast but ineffective heuristic *Pseudocost Branching* (PB). We also include the GNN model from Gasse et al. [17] run on CPU (GNN), and also several fast but less expressive models such as SVMRank from Khalil et al. [26], LambdaMART from



Figure 3: Test accuracy of the different models, with a simple e2e training protocol.

Burges [9], and ExtraTree Classifier from Geurts et al. [18] as benchmarks. For conciseness, we chose to report those last three competitor models using an optimistic aggregation scheme, by systematically choosing only the best performing method among the three (COMP). For completeness, we also
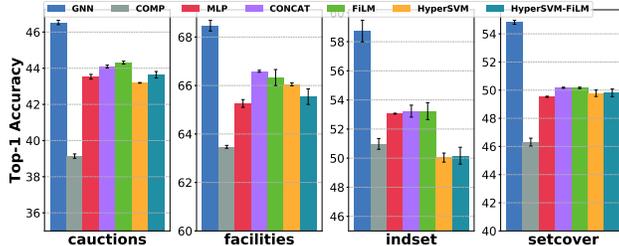
---

[4]for complete definition refer to Appendix A.3 in Achterberg [2]

report the performance of a GNN model run on a high-end GPU, although we do not consider that method as a baseline and therefore do not include it in the Wins indicator. We acknowledge that our comparison to general branching strategy like RPB is not completely fair, however, developing specialized versions of such strategies is a challenge in itself full of modeling choices that we foresee as future work.

**Model selection.** To investigate the effectiveness of different architectures, we empirically compare the performance of their end-to-end variants. Figure 3 compares the Top-1 test accuracy of models across the four problem sets. The performance of GNNs (blue), being the most expressive model, serves as an upper bound to the performance of hybrid models. All of the considered hybrid models outperform MLPs (red) across all problem sets. Additionally, we observe that FiLM (green) and CONCAT (purple) perform significantly better than other architectures. However, there is no clear winner among them. We also note that the cheapest hybrid models, HyperSVM and HyperSVM-FiLM, though better than MLPs, still do not perform as well as FiLM or CONCAT models.

**Training protocols.** In Table 2, we show the effect of different protocols discussed in section 4 on Top-1 accuracy of the FiLM models. We observe that the presence of these protocols improves the model accuracy by 0.5-0.9%, which translates to a minor yet practically useful improvement in B&B performance of the solver. Except for Combinatorial Auctions, FiLM's performance is improved by knowledge distillation, which suggests that the soft targets of pre-trained GNN yield a better generalization performance. Lastly, we launch a hyperparameter search for auxiliary objective: ED and MHE, on top of the best performing model (Top-1 accuracy), among e2e and e2e & KD models. Auxiliary tasks further help the accuracy of the hybrid models, but for some problem classes it is ED that works well while for others it is MHE. We provide the model performances on test dataset in the supplement for further reference.

**Effect of loss weighting.** We empirically investigate the effect of the different loss weighting schemes discussed in Section 4.2.4. We train a simple MLP model on our small Combinatorial Auctions instances, and measure the resulting B&B tree size on big instances. We report aggregated results in Table 3, and provide instance-level results in the supplement. We observe that the most commonly used exponential and linear schemes actually seem to degrade the performance of the learned strategy, as we believe those may be too aggressive at disregarding nodes early on in the tree.

Table 2: Test accuracy of FiLM, using different training protocols.

|  | cauctions | facilities | indset | setcover |
| --- | --- | --- | --- | --- |
| Pretrained GNN | $44.12 \pm 0.09$ | $65.78 \pm 0.06$ | $53.16 \pm 0.51$ | $50.00 \pm 0.09$ |
| e2e | $44.31 \pm 0.08$ | $66.33 \pm 0.33$ | $53.23 \pm 0.58$ | $50.16 \pm 0.05$ |
| e2e & KD | $44.10 \pm 0.09$ | $66.60 \pm 0.21$ | $53.08 \pm 0.3$ | $50.31 \pm 0.19$ |
| e2e & KD & AT | $\mathbf{44.56 \pm 0.13}$ | $\mathbf{66.85 \pm 0.28}$ | $\mathbf{53.68 \pm 0.23}$ | $\mathbf{50.37 \pm 0.03}$ |

Table 3: Effect of different sample weighting schemes on combinatorial auctions (big) instances, with a simple MLP model. $z \in [0, 1]$ is the ratio of the depth of the node and the maximum depth observed in a tree.

| Type | Weighting scheme | Nodes | Wins |
| --- | --- | --- | --- |
| Constant | $1$ | 9678 | 10/60 |
| Exponential decay | $e^{-0.5z}$ | 9793 | 10/60 |
| Linear | $(e^{-0.5} - 1) * z + 1$ | 9789 | 12/60 |
| Quadratic decay | $(e^{-0.5} - 1) * z^2 + 1$ | 9561 | 14/60 |
| Sigmoidal | $(1 + e^{-0.5})/(1 + e^{z-0.5})$ | **9534** | 14/60 |

On the other hand, both the quadratic and sigmoidal schemes result in an improvement, thus validating the idea that depth-dependent weighting can be beneficial for learning to branch. We therefore opt for a sigmoidal loss weighting scheme in our training protocol.

**Complete benchmark.** Finally, to evaluate the runtime performance of our hybrid approach, we replace SCIP's default branching strategy with our best performing model, FiLM. We observe in Table 4 that FiLM performs substantially better than all other CPU-based branching strategies. The computationally expensive FSB, our "gold standard", becomes impractical as the size of instances grows, whereas RPB remains competitive. While GNN retains its small number of nodes, it loses in running time performance on CPU. Note that we found that FiLM models for Maximum Independent Set did initially overfit on small instances, such that the performance on larger instances degraded substantially. To overcome this issue we used weight decay [34] regularization, with a validation set of 2000 observations generated using random medium instances (not used for evaluation). We report the performance of the regularized models in the supplement, and use the best performing model to report evaluation performance on medium and big instances. We also show in the supplement that the

Table 4: Performance of *branching strategies* on evaluation instances. We report geometric mean of solving times, number of times a method won (in solving time) over total finished runs, and geometric mean of number of nodes. Refer to section 5 for more details. The best performing results are in **bold**. *Models were regularized to prevent overfitting on small instances.

| | Small | | | Medium | | | Big | | |
|---|---|---|---|---|---|---|---|---|---|
| Model | Time | Wins | Nodes | Time | Wins | Nodes | Time | Wins | Nodes |
| FSB | 42.53 | 1 / 60 | 13 | 313.33 | 0 / 59 | 75 | 997.23 | 0 / 51 | 50 |
| PB | 31.35 | 4 / 60 | 139 | 177.69 | 4 / 60 | 384 | 712.45 | 3 / 56 | 309 |
| RPB | 36.86 | 1 / 60 | **23** | 213.99 | 1 / 60 | **152** | 794.80 | 2 / 54 | **99** |
| COMP | 30.37 | 3 / 60 | 120 | 172.51 | 4 / 60 | 347 | 633.42 | 6 / 57 | 294 |
| GNN | 39.18 | 0 / 60 | 112 | 209.84 | 0 / 60 | 314 | 748.85 | 0 / 54 | 286 |
| FiLM (ours) | **24.67** | **51** / 60 | 109 | **136.42** | **51** / 60 | 325 | **531.70** | **46** / 57 | 295 |
| GNN-GPU | 28.91 | – / 60 | 112 | 150.11 | – / 60 | 314 | 628.12 | – / 56 | 286 |
| | | | | Capacitated Facility Location | | | | | |
| FSB | 27.16 | 0 / 60 | 17 | 582.18 | 0 / 45 | 116 | 2700.00 | 0 / 0 | n/a |
| PB | 10.19 | 0 / 60 | 286 | 94.12 | 0 / 60 | 2451 | 2208.57 | 0 / 23 | 82 624 |
| RPB | 14.05 | 0 / 60 | **54** | 94.65 | 0 / 60 | 1129 | 1887.70 | 7 / 27 | 48 395 |
| COMP | 9.83 | 3 / 60 | 178 | 89.24 | 0 / 60 | 1474 | 2166.44 | 0 / 21 | 52 326 |
| GNN | 17.61 | 0 / 60 | 136 | 242.15 | 0 / 60 | **1013** | 2700.17 | 0 / 0 | n/a |
| FiLM (ours) | **8.73** | **57** / 60 | 147 | **63.75** | **60** / 60 | 1131 | **1843.24** | **20** / 26 | **37 777** |
| GNN-GPU | 8.26 | – / 60 | 136 | 53.56 | – / 60 | 1013 | 1535.80 | – / 36 | 31 662 |
| | | | | Set Covering | | | | | |
| FSB | 6.12 | 0 / 60 | 6 | 132.38 | 0 / 60 | 71 | 2127.35 | 0 / 28 | 318 |
| PB | 2.76 | 1 / 60 | 234 | 25.83 | 0 / 60 | 2765 | 393.60 | 0 / 59 | 13 719 |
| RPB | 4.01 | 0 / 60 | **11** | 26.36 | 0 / 60 | 714 | **210.95** | **29** / 60 | 4701 |
| COMP | 2.76 | 0 / 60 | 82 | 29.76 | 0 / 60 | 930 | 494.59 | 0 / 54 | 5613 |
| GNN | 2.73 | 1 / 60 | 71 | 22.26 | 0 / 60 | 688 | 257.99 | 6 / 60 | **3755** |
| FiLM (ours) | **2.13** | **58** / 60 | 73 | **15.71** | **60** / 60 | **686** | 217.02 | 25 / 60 | 4315 |
| GNN-GPU | 1.96 | – / 60 | 71 | 11.70 | – / 60 | 688 | 121.18 | – / 60 | 3755 |
| | | | | Combinatorial Auctions | | | | | |
| FSB | 673.43 | 0 / 53 | 47 | 1689.75 | 0 / 20 | 10 | 2700.00 | 0 / 0 | n/a |
| PB | 172.03 | 2 / 57 | 5728 | 753.95 | 0 / 45 | 1570 | 2685.23 | 0 / 1 | 38 215 |
| RPB | 59.87 | 5 / 60 | 603 | 173.17 | 11 / 60 | **205** | 1946.51 | 9 / 21 | 2461 |
| COMP | 82.22 | 1 / 58 | 847 | 383.97 | 1 / 52 | 267 | 2393.75 | 0 / 6 | 5589 |
| GNN* | **44.07** | 15 / 60 | **331** | 625.23 | 1 / 50 | 599 | 2330.95 | 0 / 10 | **687** |
| FiLM* (ours) | 52.96 | 37 / 55 | 376 | 131.45 | 47 / 54 | 264 | **1823.29** | **12** / 15 | 1201 |
| GNN-GPU* | 31.71 | – / 60 | 331 | 63.96 | – / 60 | 599 | 1158.59 | – / 27 | 685 |
| | | | | Maximum Independent Set | | | | | |

cheapest computation model of HyperSVM/HyperSVM-FiLM do not provide any runtime advantages over FiLM. We achieve up to 26% reduction on medium instances and up to 8% reduction on big instances in overall solver running time compared to the next-best branching strategy, including both learned and classical strategies.

Finally, we note that the majority of "big" problems in Set Covering and Maximum Independent Set are not solved by any of the branching strategies. Therefore, we provide a comparison of optimality gap of unsolved instances at time out in Table 5. It is evident that FiLM models are able to close a larger optimality gap than the other branching strategies.

In the absence of significant difference in KD and KD & AT model's Top-1 accuracy (see Table 2), a natural question then to ask is: what can be a practically useful choice? To answer this question, we tested the effect of each training protocol on the final B&B performance. Although the detailed discussion is in the supplement we conclude that in the absence of significant difference in the test accuracy of models between KD and AT & KD, a preference should be made for KD models to attain better generalization.

Table 5: Mean optimality gap (lower the better) of commonly unsolved "big" instances (number of such instances in brackets).

| | setcover (33) | indset (39) |
|---|---|---|
| FSB | 0.1709 | 0.0755 |
| PB | 0.0713 | 0.0298 |
| RPB | 0.0628 | 0.0252 |
| COMP | 0.0740 | 0.0252 |
| GNN | 0.1039 | 0.0341 |
| FiLM | **0.0597** | **0.0187** |

**Limitations.** We would like to point out some limitations of our work. First, given the NP-Hard nature of MILP solving, it is fairly time consuming to evaluate performance of the trained models on the instances bigger than considered for this work. One can consider the primal-dual bound gap after a time limit as an evaluation metric for the bigger instances, but this is misaligned with the solving

time objective. Second, we have used Top-1 accuracy on the test set as a proxy for the number of nodes, but there is an inherent distribution shift because of the sequential nature of B&B that leads to out-of-distribution observations. Third, generalization to larger instances is a central problem in the design of branching strategies. Several techniques discussed in this work form only a part of the solution to this problem. On the Maximum Independent Set problem, we originally noticed a poor generalization capability, which we addressed by cross-validation using a small validation set. In future work, we plan to perform an extensive study of the effect of architecture and training protocols on generalization performance. Another experiment worth conducting would be to train on larger instances than "small" problems used in the work, in order to get a better view of how and to which point the different models are able to generalize. However, not only is this a time-consuming process, but there is also an upper limit on the size of the problems on which it is reasonable to conduct experiments, simply due to hardware constraints. Finally, although we showed the efficacy of our models on a broad class of MILPs, there may be other problem classes for which our models might not result in a substantial runtime improvements.

## 6    Conclusion

As more operations research and integer programming tools start to include ML and deep learning modules, it is necessary to be mindful of the practical bottlenecks faced in that domain. To this end, we combine the expressive power of GNNs with the computational advantages of MLPs to yield novel hybrid models for learning to branch, a central component in MILP solving. We integrate various training protocols that augment the basic MLPs and help bridge the accuracy gap with more expensive models. This competitive accuracy translates into savings in time and nodes when used in MILP solving, as compared to both default, expert-designed branching strategies and expensive GNN models, thus obtaining the "best of both worlds" in terms of the time-accuracy trade-off. More broadly, our philosophy revolves around understanding the intricacies and practical constraints of MILP solving and carefully adapting deep learning techniques therein. We believe that this integrative approach is crucial to the adoption of statistical learning in exact optimization solvers.

### Broader Impact

This paper establishes a bridge between the work done in ML in the last years on learning to branch and the traditional MILP solvers used to routinely solve thousands of optimization applications in energy, telecommunications, logistics, biology, just to mention a few.

The MILP solvers are executed on CPU-only machines and the technical challenge of using GPU-based algorithmic techniques to hybridize them had been neglected thus far. Admittedly, such a challenge was not urgent when the learning to branch literature was in its early stages. That situation has changed drastically with the GNN implementation in [17], the first approach to show significant benefit with respect to the default version of a state-of-the-art MILP solver like SCIP. For this reason, the current paper comes at the due time for the literature in the field and addresses the challenge, for the first time, in a sophisticated, yet relatively simple way.

Thus, our work provides the first viable way for commercial and noncommercial MILP solver developers to implement and integrate a ML-based "learning to branch" framework and for hundreds of thousands of users and practitioners to use it. In an even broader sense, the fact that we were able to approximate the performance of GPU-based models with a sophisticated integration of CPU-based techniques is consistent with, for example, Hinton et al. [25], and widens the space of problems to which ML techniques can be successfully applied.

Yu-Tung Hui, Tristan Deleu, Maksym Korablyov, and Alex Lamb for enlightening discussions on deep learning and integer programming.

## References

[1] Jawad Abrache, Teodor Gabriel Crainic, Michel Gendreau, and Monia Rekik. Combinatorial auctions. *Annals of Operations Research*, 153(1):131–164, 2007.

[2] Tobias Achterberg. *Constraint Integer Programming*. Doctoral thesis, Technische Universität Berlin, Fakultät II - Mathematik und Naturwissenschaften, Berlin, 2007. URL http://dx.doi.org/10.14279/depositonce-1634.

[3] Tobias Achterberg and Timo Berthold. Hybrid branching. In *Integration of AI and OR Techniques in Constraint Programming for Combinatorial Optimization Problems*, 2009.

[4] Tobias Achterberg, Thorsten Koch, and Alexander Martin. Branching rules revisited. *Operations Research Letters*, 33(1):42 – 54, 2005. ISSN 0167-6377. doi: https://doi.org/10.1016/j.orl.2004.04.002. URL http://www.sciencedirect.com/science/article/pii/S0167637704000501.

[5] Alejandro Marcos Alvarez, Quentin Louveaux, and Louis Wehenkel. A machine learning-based approximation of strong branching. *INFORMS Journal on Computing*, 29(1):185–195, 2017.

[6] David Applegate, Robert Bixby, Vašek Chvátal, and William Cook. Finding cuts in the TSP. Technical report, DIMACS, 1995.

[7] Yoshua Bengio, Andrea Lodi, and Antoine Prouvost. Machine learning for combinatorial optimization: a methodological tour d'horizon. *arXiv:1811.06128*, 2018.

[8] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016.

[9] Christopher JC Burges. From ranknet to lambdarank to lambdamart: An overview. *Learning*, 11(23-581):81, 2010.

[10] Bhuwan Dhingra, Hanxiao Liu, Zhilin Yang, William Cohen, and Ruslan Salakhutdinov. Gated-attention readers for text comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1832–1846, 2017.

[11] Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur. A learned representation for artistic style. *arXiv preprint arXiv:1610.07629*, 2016.

[12] Vincent Dumoulin, Ethan Perez, Nathan Schucher, Florian Strub, Harm de Vries, Aaron Courville, and Yoshua Bengio. Feature-wise transformations. *Distill*, 3(7):e11, 2018.

[13] DG Elson. Site location via mixed-integer programming. *Journal of the Operational Research Society*, 23 (1):31–43, 1972.

[14] Frank Finn. Integer programming, linear programming and capital budgeting. *Abacus*, 9:180 – 192, 07 2005. doi: 10.1111/j.1467-6281.1973.tb00186.x.

[15] Christodoulos A Floudas and Xiaoxia Lin. Mixed integer linear programming in process scheduling: Modeling, algorithms, and applications. *Annals of Operations Research*, 139(1):131–162, 2005.

[16] J.J. Thomson F.R.S. XXIV. on the structure of the atom: an investigation of the stability and periods of oscillation of a number of corpuscles arranged at equal intervals around the circumference of a circle; with application of the results to the theory of atomic structure. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 7(39):237–265, 1904. doi: 10.1080/14786440409463107. URL https://doi.org/10.1080/14786440409463107.

[17] Maxime Gasse, Didier Chételat, Nicola Ferroni, Laurent Charlin, and Andrea Lodi. Exact combinatorial optimization with graph convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 15554–15566, 2019.

[18] Pierre Geurts, Damien Ernst, and Louis Wehenkel. Extremely randomized trees. *Machine learning*, 63(1):3–42, 2006.

[19] Tobias Glasmachers. Limits of end-to-end learning. *arXiv preprint arXiv:1704.08305*, 2017.

[20] Ambros Gleixner, Michael Bastubbe, Leon Eifler, Tristan Gally, Gerald Gamrath, Robert Lion Gottwald, Gregor Hendel, Christopher Hojny, Thorsten Koch, Marco E. Lübbecke, Stephen J. Maher, Matthias Miltenberger, Benjamin Müller, Marc E. Pfetsch, Christian Puchert, Daniel Rehfeldt, Franziska Schlösser, Christoph Schubert, Felipe Serrano, Yuji Shinano, Jan Merlin Viernickel, Matthias Walter, Fabian Wegscheider, Jonas T. Witt, and Jakob Witzig. The SCIP Optimization Suite 6.0. Technical report, Optimization Online, July 2018. URL `http://www.optimization-online.org/DB_HTML/2018/07/6692.html`.

[21] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. `http://www.deeplearningbook.org`.

[22] David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016.

[23] He He, Hal III Daumé, and Jason Eisner. Learning to search in branch-and-bound algorithms. In *Advances in Neural Information Processing Systems 27*, pages 3293–3301, 2014.

[24] Marti A. Hearst. Support vector machines. *IEEE Intelligent Systems*, 13(4):18–28, July 1998. ISSN 1541-1672. doi: 10.1109/5254.708428. URL `https://doi.org/10.1109/5254.708428`.

[25] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.

[26] Elias B. Khalil, Pierre Le Bodic, Le Song, George Nemhauser, and Bistra Dilkina. Learning to branch in mixed integer programming. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, pages 724–731, 2016.

[27] Taesup Kim, Inchul Song, and Yoshua Bengio. Dynamic layer normalization for adaptive neural acoustic modeling in speech recognition. *Proc. Interspeech 2017*, pages 2411–2415, 2017.

[28] Solomon Kullback and Richard A Leibler. On information and sufficiency. *The annals of mathematical statistics*, 22(1):79–86, 1951.

[29] A. H. Land and A. G. Doig. An automatic method of solving discrete programming problems. *Econometrica*, 28(3):pp. 497–520, 1960.

[30] Lukas Liebel and Marco Körner. Auxiliary tasks in multi-task learning. *arXiv preprint arXiv:1805.06334*, 2018.

[31] Jeff Linderoth and Martin Savelsbergh. A computational study of search strategies for mixed integer programming. *INFORMS Journal on Computing*, 11:173–187, 05 1999. doi: 10.1287/ijoc.11.2.173.

[32] Weiyang Liu, Rongmei Lin, Zhen Liu, Lixin Liu, Zhiding Yu, Bo Dai, and Le Song. Learning towards minimum hyperspherical energy. In *Advances in neural information processing systems*, pages 6222–6233, 2018.

[33] Andrea Lodi and Giulia Zarpellon. On learning and branching: a survey. *TOP*, 25:207–236, 2017.

[34] Andrew Y Ng. Feature selection, l1 vs. l2 regularization, and rotational invariance. In *Proceedings of the twenty-first international conference on Machine learning*, page 78, 2004.

[35] Christos H. Papadimitriou and Kenneth Steiglitz. *Combinatorial Optimization: Algorithms and Complexity*. Prentice-Hall, Inc., USA, 1982. ISBN 0131524623.

[36] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. FiLM: Visual reasoning with a general conditioning layer. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

[37] Rupesh Kumar Srivastava, Klaus Greff, and Jürgen Schmidhuber. Highway networks. *arXiv preprint arXiv:1505.00387*, 2015.

[38] Laurence A. Wolsey. *Integer Programming*. Wiley-Blackwell, 1988.

[39] Giulia Zarpellon, Jason Jo, Andrea Lodi, and Yoshua Bengio. Parameterizing branch-and-bound search trees to learn branching policies. *arXiv preprint arXiv:2002.05120*, 2020.

# Supplement: Hybrid Models for Learning to Branch

**Prateek Gupta**[*]
University of Oxford
The Alan Turing Institute

**Maxime Gasse**
Mila, Polytechnique Montréal

**Elias B. Khalil**
University of Toronto

**M. Pawan Kumar**
University of Oxford

**Andrea Lodi**
CERC, Polytechnique Montréal

**Yoshua Bengio**
Mila, Université de Montréal

## 1 Inefficiency in using GNNs for solving MILPs in parallel

In this section, we argue that the GNN architecture looses its advantages in the face of solving multiple MILPs at the same time. In the applications like multi-objective optimization [4], where multiple MILPs are solved in parallel, a GNN for each MILP needs to be initialized on the GPU because of the *sequentially asynchronous* nature of solving MILPs. Not only is there a limit to the number of such GNNs that can fit on a single GPU because of memory constraints, but also several GNNs on a single GPU results in an inefficient GPU utilization.

One can, for instance, try to time multiple MILPs such that there is a need for a single forward evaluation on a GPU, but, in our knowledge, it has not been done and it results in frequent interruptions in the solving procedure. An alternative, much simpler, method is to *pack* multiple GNNs on a single GPU such that each GNN is dedicated to solving one MILP. For example, we were able to put 25 GNNs on Tesla V100 32 GB GPU. Figure 1 shows the inefficient utilization of GPUs when multiple GNNs are packed on a single GPU.
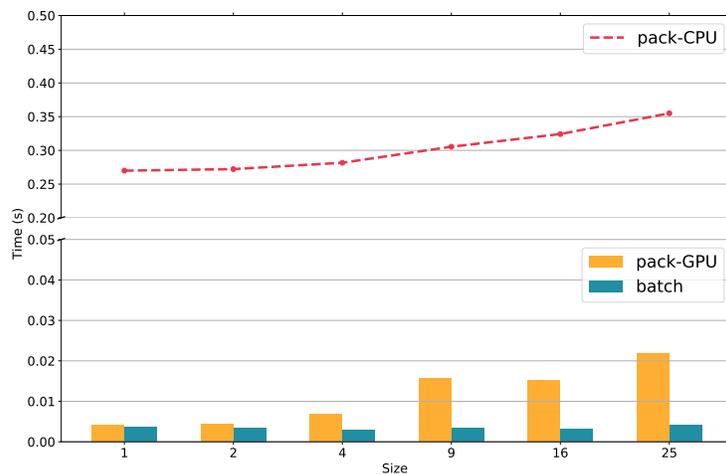


Figure 1: *Packing* several GNNs together on a GPU keeps it underutilized. "Size" is the number of inputs batched together (unrealistic scenario) or number of inputs simultaneously put on a GPU separately.

---

[*]The work was done during an internship at Mila and CERC. Correspondence to: <pgupta@robots.ox.ac.uk>

# 2 Input Features

We use the features that were used by Gasse et al. [1] and Khalil et al. [3]. The list of features in **G** are described in the Table 1. There are a total of 92 features that we use for **X** as described in the Table 2. We follow the preprocessing procedure as described in Gasse et al. [1].

Table 1: Features of **G**. It is constructed as a variable-constraint bipartite graph, where the vertices of this graph have features as described here. These features are same as used in Gasse et al. [1].

|  | Type | Description |
|---|---|---|
| **Variable Features (Total -13)** | | |
| Variable type (1) | Categorical | Allowed values - Binary, Integer, Continuous, Implied Integer |
| Normalized coefficient (1) | Real | Objective coefficient of the variable normalized by the euclidean norm of its coefficients in the constraints |
| Specified bounds (2) | Binary | Does the variable has a lower bound (upper bound)? |
| Solution bounds (2) | Binary | If the variable is currently at its lower bound (upper bound)? |
| Solution bractionality (1) | Real $\in [0, 1)$ | Fractional part of the variable i.e. $x - \lfloor x \rfloor$, where $x$ is the value of the decision in variable in the current LP solution |
| Basis (1) | Categorical | One of 4 classes - Lower (variable is at the lower bound), Basic (variable has a value between the bounds), Upper (variable is at the upper bound), Zero (rare case) |
| Reduced cost (1) | Real | Amount by which objective coefficient of the variable should decrease so that the variable assumes a positive value in the LP solution |
| Age (1) | Real | Number of LP iterations since the last time the variable was basic normalized by total number of LP iterations |
| Solution value (1) | Real | Value of the variable in the current LP solution |
| Primal value (1) | Real | Value of the variable in the current best primal solution |
| Average primal value (1) | Real | Average value of the variable in all of the previously observed feasible primal solutions |
| **Constraint Features (Total - 5)** | | |
| Cossine similarity (1) | Real | Cossine of angle between the vector represented by objective coefficients and the coefficients of this constraint |
| Bias (1) | Real | Right hand side of the constraint normalized by the euclidean norm of the row coefficients |
| Age (1) | Real | Number of iterations since the last time the constraint was active normalized by total number of LP iterations |
| Normalized dual value (1) | Real | Value of dual variable corresponding to the constraint normalized by the product of norms of the row coefficients and the objective coefficients |
| Bounds (1) | Binary | If the constraint is currently at its bounds? |
| **Edge Features (Total - 1)** | | |
| Normalized coefficient (1) | Real | Coefficient of the variable normalized by the norm of the coefficients of all the variables in the constraint |

Table 2: Features in **X**, an input to MLP.

|  | count |
|---|---|
| Variable features from **G** | 13 |
| Variable features from Khalil et al. [3] | 72 |

# 3 Preliminary results on running times of various architectures

In order to have a rough idea of relative improvement in runtimes across various architectures, we considered 20 instances in each difficulty level for each problem class and use the solver to obtain observations at each node. Different functions are then used to evaluate decisions at each node, and the total time taken for all the function evaluations across the nodes in an instance is observed. Note that the quality of decision is not important here; we are only interested in time taken per decision.

Specifically, we considered 5 forms of architectures as listed in Table 3. To cover the entire spectrum of strong inductive biases, we constructed an attention mechanism as explained in 4. At the same time, to cover the spectrum of cheaper architectures we consider simple dot product as the form of predictor. Finally, we consider a scenario where we only consider MLPs as predictors everywhere in the B&B tree.

Table 3: Different architectures considered for preliminray runtime comparison.

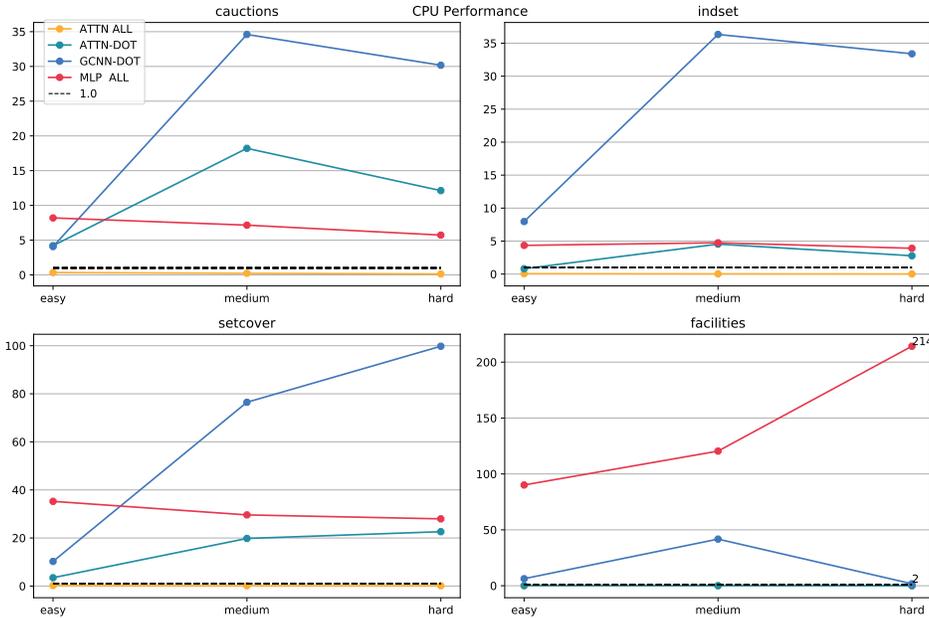| Type | Description |
|---|---|
| GNN ALL | Use a Graph Convolution Neural Network (GNN) [1] at all *tree nodes* |
| ATTN ALL | Attention mechanism (section 4) at all nodes |
| GNN DOT | Use GNN at the *root nodes* and a dot product with $\mathbf{X}$ to compute scores at every node |
| ATTN DOT | Use attention mechanism at the root node and dot product with $\mathbf{X}$ to compute scores at every node |
| MLP ALL | Use a 3-layer multilinear perceptron at all the nodes |



Figure 2: Relative time performance of various methods with respect to GNN ALL on CPU (average values). A value of 10 implies that the method is 10 times faster (on arithmetic average) than GNN ALL implying that one can afford to perform 10 times worse than GNN ALL in iterative performance when using CPUs. Note that this is just a rough estimation.

Figure 2 shows the relative performance (rel. to GNN) of various deep learning architectures across 4 sets of problems. It is evident that MLP ALL and GNN DOT are favored across the problem sets. This observation inspired the range of *hybrid architectures* that we explored in the paper. We also observe that a superior inductive bias like that of attention mechanism and transformers [7] has a better runtime performance on GPUs as illustrated in Figure 3, which we suspect is massive parallelization employed in the computations of attention. However, their performance on CPUs is not better than GNNs.
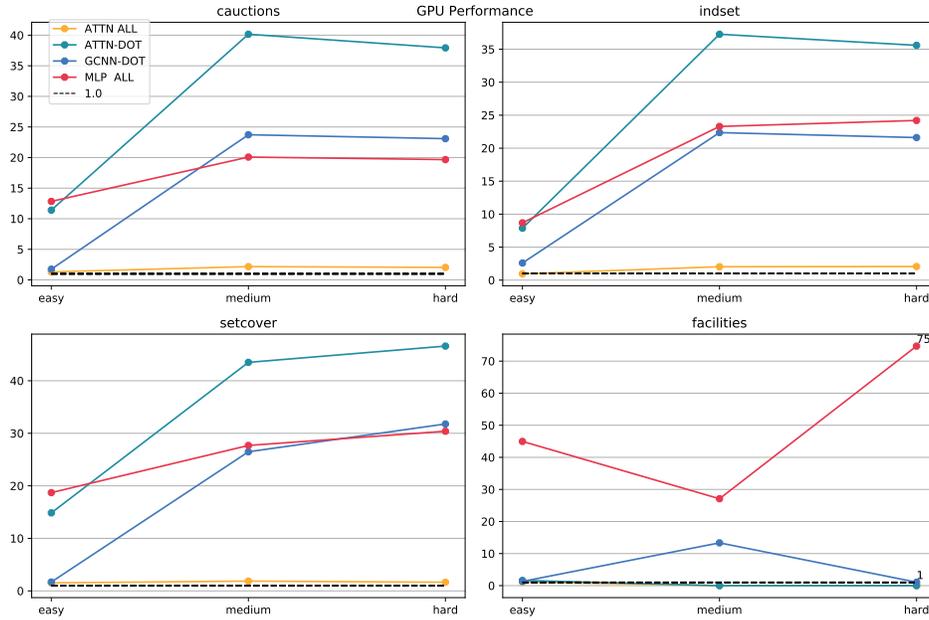
Figure 3: Relative time performance of various methods with respect to $GNNALL$ on GPU (average values). A value of $10$ implies that the method is $10$ times faster (on arithmetic average) than $GNNALL$ implying that one can afford to perform $10$ times worse than $GNNALL$ in iterative performance when using CPUs. Note that this is just a rough estimation.

## 4 Attention Mechanism for MILPs

To cover the entire spectrum of computational complexity and expressivity of inductive bias, we implemented a transformer [7] as an architecture to replace GNNs. Specifically, we let the variables (constraints) attend to all other variables (constraints) via multi-headed self-attention mechanism. Finally, a modulated attention mechanism between *variable representations* as queries and *constraint representation* as keys outputs final *variable representations*, which is passed through the softmax layer for classification objective. Here, we use modulation scheme as explained in Shaw et al. [6]. Precisely, an edge in variable-constraint graph is used to increment the attention score with a learnable scalar value.
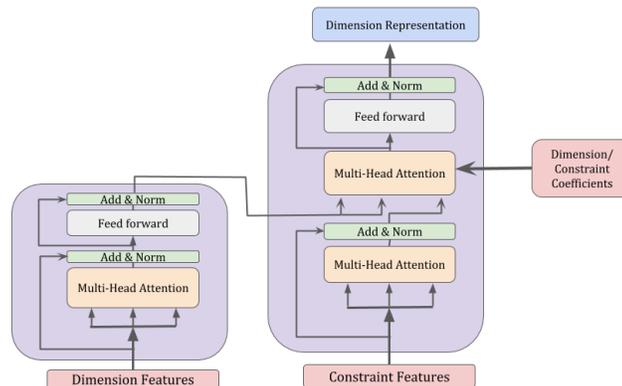


Figure 4: Multi-Head Attention mechanism where a variable or constraint can attend to all other variables or constraints. Finally,variables are used as a query to attend to constraints, where the attention is modulated through variable-constraint features. Modulation of attention scores follow Shaw et al. [6]

4

# 5 Data Generation & Training Specifications

For our experiments, we used 10,000 training instances for the training dataset and collected 150,000 observations from the tree nodes of those instances. In a similar manner, we generated 20,000 instances each for the validation and testing set, resulting in 30,000 observations each for validation and testing respectively.

We held all the training parameters fixed across the models. Specifically, we used a learning rate of $1e^{-3}$, training batch size of 32, a learning rate scheduler to reduce learning rate by $0.2$ if there is no improvement in the validation loss for $15$ epochs, and an early stopping criterion of $30$ epochs. Our epoch consisted of 10K training examples and 2K validation samples. We used a 3 layered MLP with 256 hidden units in each layer, while we used GNN model with an embedding size of 64 units. We used this configuration across all the models discussed in the main paper. Due to the large size of instances in capacitated facility location, we used a learning rate of 0.005, early stopping criterion of 20 epochs with a patience of 10 epochs.

Further, for knowledge distillation we used $T = 2$ (temperature) and $\alpha = 0.9$ (convex mixing of soft and hard objectives). These are the recommended settings in Hinton et al. [2]. We did a hyperparameter search for $\beta = \{0.01, 0.001, 0.0001\}$ for ED and MHE. Following are the values for $\beta$ that resulted in the best performing models.

Table 4: Best AT models

|  | AT |
| --- | --- |
| Capacitated Facility Location | MHE |
| Combinatorial Auctions | MHE |
| Set Covering | ED |
| Maximum Independent Set | MHE |

We implemented all the models using PyTorch [5], and ran all the CPU evaluations on an Intel(R) Xeon(R) CPU E5-2650 v4 @ 2.20GHz. GPU evaluations for GNN were performed on NVIDIA-TITAN Xp GPU card with CUDA 10.1.

# 6 Depth-dependent loss weighting scheme

In Figure 5 we plot the sorted ratio of number of nodes to the minimum number of nodes observed for an instance across all the weighting schemes. Here we attempt to breakdown the performance of different branching strategies learned using different loss-weighting schemes. We observe that the *sigmoidal* scheme achieves the best performance.
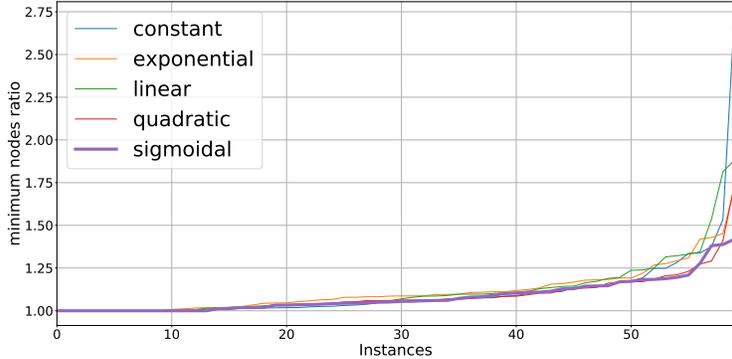


Figure 5: Performance of different weighing schemes across "big" instances. "Sigmoidal" evolves slowest among all.

# 7 Model Results

The table below is a list of all the architectures along with their performance on the test sets. Please note that training with auxiliary task is not possible with HyperSVM type of architecture so their results are not reported in the table.

Table 5: Top-1 accuracy of various models.

|  |  | Combinatorial Auctions | Capacitated Facility Location | Set Cover | Maximum Independent Set |
|---|---|---|---|---|---|
| Expert | GCNN | $46.53 \pm 0.12$ | $68.47 \pm 0.22$ | $54.82 \pm 0.14$ | $58.73 \pm 0.54$ |
| Existing methods | Extratrees | $37.97 \pm 0.1$ | $60.06 \pm 0.18$ | $42.66 \pm 0.22$ | $19.35 \pm 0.55$ |
|  | LMART | $38.7 \pm 0.16$ | $63.28 \pm 0.06$ | $46.31 \pm 0.28$ | $50.98 \pm 0.37$ |
|  | SVMRank | $39.14 \pm 0.12$ | $63.47 \pm 0.06$ | $46.23 \pm 0.11$ | $49.29 \pm 0.23$ |
| Simple Network | MLP | $43.53 \pm 0.13$ | $65.26 \pm 0.16$ | $49.53 \pm 0.04$ | $53.06 \pm 0.03$ |
| Concatenate | CONCAT (Pre) | $44.16 \pm 0.03$ | $65.5 \pm 0.1$ | $49.98 \pm 0.18$ | $52.85 +- 0.34$ |
|  | CONCAT (e2e) | $44.09 \pm 0.08$ | $66.59 \pm 0.04$ | $50.17 \pm 0.04$ | $53.23 \pm 0.41$ |
|  | CONCAT (e2e & KD) | $44.19 \pm 0.05$ | $66.53 \pm 0.13$ | $50.11 \pm 0.09$ | $53.66 +- 0.32$ |
| Modulate | FiLM (Pre) | $44.12 \pm 0.09$ | $65.78 \pm 0.06$ | $50.0 \pm 0.09$ | $53.16 \pm 0.51$ |
|  | FiLM (e2e) | $44.31 \pm 0.08$ | $66.33 \pm 0.33$ | $50.16 \pm 0.05$ | $53.23 \pm 0.58$ |
|  | FiLM (e2e & KD) | $44.1 \pm 0.09$ | $66.6 \pm 0.21$ | $50.31 \pm 0.19$ | $53.08 \pm 0.3$ |
| HyperSVM | HyperSVM (e2e) | $43.19 \pm 0.02$ | $66.05 \pm 0.06$ | $49.78 \pm 0.23$ | $50.04 \pm 0.31$ |
|  | HyperSVM (e2e & KD) | $42.55 \pm 0.03$ | $66.07 \pm 0.05$ | $49.53 \pm 0.13$ | $49.34 +- 0.43$ |
| HyperSVM-FiLM | HyperSVM-FiLM (e2e) | $43.64 \pm 0.18$ | $65.54 \pm 0.32$ | $49.81 \pm 0.27$ | $50.17 \pm 0.58$ |
|  | HyperSVM-FiLM (e2e & KD) | $43.28 \pm 0.48$ | $65.52 \pm 0.34$ | $49.73 \pm 0.05$ | $49.73 +- 0.39$ |
|  | FiLM (e2e & KD & AT) | $\mathbf{44.56 \pm 0.13}$ | $\mathbf{66.85 \pm 0.28}$ | $\mathbf{50.37 \pm 0.03}$ | $\mathbf{53.68 \pm 0.23}$ |

# 8 Overfitting in maximum independent set

Table 6 shows that the FiLM models overfit on small instances of maximum independent set. To address this problem, we used a mini dataset of 2000 observations obtained by running data collection on medium instances of maximum independent set. Further, we regularized the FiLM parameters of the FiLM model to yield much simpler models based on the performance on this mini-dataset, which is not too expensive to obtain owing to the size of observations. The results of the regularized models are in Table 7 For a fair comparison, we regularized GNN and used the best performing model to report evaluation results in the main paper.

Table 6: FiLM models for maximum independent set overfits on small instances

|  | small | | | medium | | | big | | |
|---|---|---|---|---|---|---|---|---|---|
|  | Time | Wins | Nodes | Time | Wins | Nodes | Time | Wins | Nodes |
| FiLM | 52.96 | **39**/ 55 | 492 | 1515.19 | 9/ 17 | 2804 | 2700.02 | 0/ 0 | nan |
| GNN-CPU | **44.07** | 21/ 60 | **432** | **371.81** | **36**/ 40 | **558** | **1981.43** | **10**/ 10 | **6334** |
| GNN-GPU | 31.70 | 0/ 59 | 432 | 264.01 | 0/ 43 | 558 | 1772.12 | 0/ 13 | 6313 |

Table 7: Top-1 accuracy of regularized models on 2000 observations from medium random instances of maximum independent set.

| weight decay | FiLM | GNN |
|---|---|---|
| 1.0 | $55.15 \pm 0.07$ | $31.6 \pm 6.63$ |
| 0.1 | $\mathbf{56.13 \pm 0.32}$ | $\mathbf{37.23 \pm 0.92}$ |
| 0.01 | $53.25 \pm 0.95$ | $26.8 \pm 16.71$ |
| 0.0 | $19.18 \pm 4.24$ | $34.08 \pm 4.8$ |

# 9 Performance of HyperSVM architectures

HyperSVM architectures are the cheapest in computation. In this section we ask if we are willing to lose machine learning accuracy by 1-2% in HyperSVM architecture, can we still get faster running times with HyperSVM type of architectures? Table 8 compares solver performance by using HyperSVM architecture and FiLM architecture. We observe that HyperSVM types of architectures loose their ability to generalize on larger scale instances.

Table 8: FiLM architecture generalizes to larger instances better than the HyperSVM types of architectures. We compare the performance of the best FiLM architecture from the Table 5 with the best HyperSVM architecture in the same table.

|  | small | | | medium | | | big | | |
|---|---|---|---|---|---|---|---|---|---|
|  | Time | Wins | Nodes | Time | Wins | Nodes | Time | Wins | Nodes |
| FiLM | **24.67** | **53**/ 60 | **109** | **136.42** | 51/ 60 | **336** | **531.70** | **46**/ 57 | **345** |
| HyperSVM | 27.26 | 7/ 60 | 110 | 158.97 | 9/ 60 | 345 | 614.36 | 11/ 57 | 346 |
| MLP | 27.61 | 4/ 60 | 114 | 156.30 | 11/ 60 | 347 | 595.31 | 9/ 56 | 334 |

(a) Capacitated Facility Location

|  | small | | | medium | | | big | | |
|---|---|---|---|---|---|---|---|---|---|
|  | Time | Wins | Nodes | Time | Wins | Nodes | Time | Wins | Nodes |
| FiLM | **8.73** | **59**/ 60 | **147** | **63.75** | **60**/ 60 | **2169** | **1843.24** | **24**/ 26 | **38530** |
| HyperSVM-FiLM | 9.73 | 1/ 60 | 148 | 72.53 | 0/ 60 | 2217 | 2061.56 | 1/ 22 | 47277 |
| MLP | 9.98 | 0/ 60 | 157 | 77.48 | 0/ 60 | 2299 | 1984.26 | 1/ 24 | 40188 |

(b) Set Cover

# 10 Scaling to twice the size of Big instances

In this section we investigate the generalization power of FiLM models trained on dataset obtained from small instances. Specifically, we generate 20 random instances of size double that of big instances, and use the trained models of FiLM and GNN to compare their performance against RPB. We use the time limit of 7200 seconds, i.e. 2 hours to account for longer solving running times as one scales out to bigger instances. We observe that the power to generalize is largely dependent on the problem family. For example, FiLM models can still outperform other strategies on scaling out on capacitated facility location problems. However, we found that RPB remains competitive on setcover problems. We note that this shortcoming of the FiLM models can be overcome via larger size of hidden layers and training on slightly larger instances than what has been used in the main paper. Table 9 shows these results.

Table 9: Performance of branching strategies on twice the size of biggest instances ($2 \times$ Big) considered in the main paper. 20 "Bigger" instances were solved using 3 seeds each resulting in a total of 60 runs. We see that FiLM models still remain competitive, and it is highly dependent on the family of problem.

|  | facilities | | | setcover | | |
|---|---|---|---|---|---|---|
| Model | Time | Wins | Nodes | Time | Wins | Nodes |
| RPB | 7200.14 | 0 / 0 | n/a | **6346.17** | **9** / 12 | 135 132 |
| GNN | 7111.83 | 1 / 5 | n/a | 7200.25 | 0 / 0 | n/a |
| FiLM (ours) | **7052.27** | **4** / 4 | n/a | 6508.78 | 3 / 10 | **107 187** |
| GNN-GPU | 6625.55 | – / 13 | n/a | 6008.14 | – / 12 | 93 909 |

# 11 Effect of different training protocols on B&B performance

Table 10 shows the performance of different training protocols on B&B performance. Although the models trained with auxiliary tasks (AT) and knowledge distillation (KD) are a clear winner up until problem sets of size medium, there is a tie between the models trained with KD and those trained with KD & AT. However, it is worth noting that the difference in B&B performance across different training protocols is not huge, and such performance evaluation can be read from the test accuracy of trained models. For example, as evident in Table 5, difference in test accuracy of FiLM (e2e & KD)

Table 10: Effect of training protocols on the performance of *branching strategies*. We report geometric mean of solving times, number of times a method won (in solving time) over total finished runs, and geometric mean of number of nodes. Refer to section 5 (main) for more details. The best performing results are in **bold**. *Models were regularized to prevent overfitting on small instances.

| Model | Small | | | Medium | | | Big | | |
|---|---|---|---|---|---|---|---|---|---|
| | Time | Wins | Nodes | Time | Wins | Nodes | Time | Wins | Nodes |
| e2e | 25.05 | 22 / 60 | 114 | 154.12 | 4 / 60 | 340 | 550.03 | 15 / 57 | 339 |
| e2e + KD | 28.00 | 2 / 60 | 111 | 143.12 | 18 / 60 | 343 | **507.50** | **26** / 57 | **325** |
| e2e + KD + AT | **24.67** | **36** / 60 | **109** | **136.42** | **38** / 60 | **336** | 531.70 | 16 / 57 | 345 |
| | Capacitated Facility Location | | | | | | | | |
| e2e | 9.70 | 1 / 60 | 152 | 71.18 | 1 / 60 | 2186 | 1869.57 | 4 / 25 | **40 341** |
| e2e + KD | 9.81 | 0 / 60 | **146** | 70.88 | 4 / 60 | 2173 | **1842.29** | **15** / 25 | 40 437 |
| e2e + KD + AT | **8.73** | **59** / 60 | 147 | **63.75** | **55** / 60 | **2169** | 1843.24 | 8 / 26 | 40 881 |
| | Set Covering | | | | | | | | |
| e2e | 2.45 | 1 / 60 | **72** | 17.59 | 5 / 60 | 702 | 225.88 | 17 / 60 | 8939 |
| e2e + KD | 2.35 | 0 / 60 | 73 | 17.59 | 2 / 60 | 720 | 241.95 | 7 / 59 | 8846 |
| e2e + KD + AT | **2.13** | **59** / 60 | 73 | **15.71** | **53** / 60 | **686** | **217.02** | **36** / 60 | **8711** |
| | Combinatorial Auctions | | | | | | | | |
| e2e | 205.63 | 2 / 54 | 559 | 1103.47 | 1 / 29 | 1137 | 2457.94 | 1 / 4 | **2268** |
| e2e + KD | 333.52 | 1 / 52 | 486 | 926.12 | 1 / 41 | 604 | 2503.65 | 0 / 7 | 1953 |
| e2e + KD + AT | **52.96** | **54** / 55 | **410** | **131.45** | **54** / 54 | **331** | **1823.29** | **14** / 15 | 3049 |
| | Maximum Independent Set* | | | | | | | | |

and FiLM (e2e & KD & AT) is not significant for problem sets - Capacitated Facility Location and Set Covering, but its significant enough for problem sets - Combinatorial Auctions and Maximum Independent Set. **Given that the inference cost is independent of training protocols, to attain better generalization performance, we recommend FiLM (e2e & KD & AT) only when these models have significantly better accuracy than FiLM (e2e & KD).**

## References

[1] Maxime Gasse, Didier Chételat, Nicola Ferroni, Laurent Charlin, and Andrea Lodi. Exact combinatorial optimization with graph convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 15554–15566, 2019.

[2] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.

[3] Elias B. Khalil, Pierre Le Bodic, Le Song, George Nemhauser, and Bistra Dilkina. Learning to branch in mixed integer programming. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, pages 724–731, 2016.

[4] Gokhan Kirlik and Serpil Sayın. A new algorithm for generating all nondominated solutions of multiobjective discrete optimization problems. *European Journal of Operational Research*, 232(3):479–488, 2014.

[5] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.

[6] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. *arXiv preprint arXiv:1803.02155*, 2018.

[7] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.