GAA-PTRNET: GRAPH ATTENTION AGGREGATION-BASED POINTER NETWORK FOR ONE-SHOT DAG SCHEDULING

Anonymous authorsPaper under double-blind review

000

001

003

004

006

008 009 010

011 012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

032

034

037 038

039

040

041

042

043

044

046

047

048

049

051

052

ABSTRACT

Optimizing Directed Acyclic Graph (DAG) workflow makespan by scheduling techniques is a critical issue in the high performance computing area. Many studies in recent years combined Pointer Network (PtrNet) with reinforcement learning (RL) to schedule DAGs by generating DAG task priorities in a sequenceto-sequence manner. However, these PtrNet-based scheduling methods need to repeatedly compute the decoder's hidden state or context embeddings according to the recent local decisions, which leads to limited capability of exploiting the DAG global topological structure, high computation complexity and inability to achieve one-shot scheduling. To address these issues, we propose GAA-PtrNet, a novel PtrNet based on graph attention aggregation (GAA) for one-shot DAG workflow scheduling. In GAA-PtrNet, we compute the pair-wise graph attention scores among nodes in one-shot, then directly aggregate these scores to obtain the probability of selecting candidate nodes. Consequently, the explicit decoder or context embedding structure in PtrNet is omitted in our GAA-PtrNet, and the network takes only one shot forward propagation to infer a solution for a whole DAG scheduling problem, significantly reducing the computation complexity. Additionally, to train GAA-PtrNet, we design a training strategy based on policy gradient RL with dense reward signal and demonstration learning. To our knowledge, GAA-PtrNet is the first network model to achieve PtrNet-based one-shot DAG scheduling. GAA-PtrNet can better handle with DAG workflow structures, providing high quality DAG scheduling solutions. The experimental results show that the proposed method is superior in terms of objective and runs about 10 times faster when compared to previous PtrNet-based methods, and also performs better than other learning-based DAG scheduling methods.

1 Introduction

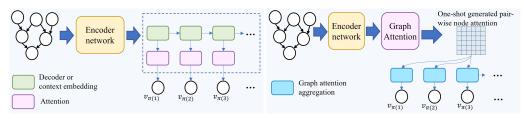
The Directed Acyclic Graph (DAG) scheduling problem arises in the high performance computing field (Hosseini Shirvani, 2024). It is a class of NP-hard Combinatorial Optimization Problems (COP), involving large and complex DAG workflow, and homogeneous or heterogeneous computation resources. In this domain, DAG is used to represent the parallel and sequential relationships among computation tasks. These tasks modeled as the nodes in the DAG, and the directed edges in the DAG represent the precedence constraints among the tasks. The goal is to achieve the best performance by determining an optimal node execution or priority order and allocating then to computation resources. In recent years, machine learning, especially reinforcement learning (RL) technique, has already shown promise in DAG scheduling (Gu et al., 2025). A common network architecture for RL-based DAG scheduling consists of a Graph Neural Network (GNN) encoder to extract workflows' structural information, and a policy network to output scheduling decisions (Mao et al., 2019) (Zhou et al., 2022) (Song et al., 2023) (Yu et al., 2023) (Dong et al., 2023).

Since DAG scheduling problem can be interpreted as a ordering problem over the problem components, Pointer Network (PtrNet) (Vinyals et al., 2015), as a sequence learning method, has shown its distinct advantages. Kintsakis et al. (2019) first introduces PtrNet into DAG workflow scheduling, which follows the foundational structure proposed by Bello et al. (2016). It encodes the task features with a long short-term memory (LSTM) encoder, and feeds the feature embeddings to an LSTM de-

coder. At each decision step, it computes the additive attention scores the current candidate actions according to the the decoder's hidden state, and the task with the highest attention score (i.e., the pointer) is selected for the following scheduling, thereby this method constructs the entire task priority sequence incrementally. Such similar PtrNet-based scheduling method is adopted by Dong et al. (2021) Zhao et al. (2022) Chen & Wang (2024) and Li et al. (2022) for DAG scheduling. In contrast, Lee et al. (2020) Lee et al. (2021) Shi & Yu (2023) and Wang et al. (2023), follows the improved PtrNet proposed by Deudon et al. (2018) and Kool et al. (2018). These works abandon calculating attention scores from LSTM decoder's hidden state. Instead, they take the context embedding of the current environment state and recent decisions to compute attention, partly addressing the limitations of traditional LSTM-based PtrNet scheduler for graphed structure problems.

However, these PtrNet-based scheduling methods still suffer from limited capability of exploiting the global topological structure of DAGs, high time complexity and inability to achieve one-shot scheduling. As shown in Fig 1a, to construct a complete solution, PtrNet requires to repeatedly calculate the decoder hidden states or the context embeddings according to the observed environment state and recent decisions, and further obtain the attention scores of the candidate decisions according to these hidden states or embeddings. Consequently, their performance are often limited by their reliance on local information from recent decisions, which fail to capture long-range dependencies and the global topological structure inherent in the DAG. The repeated calculation of decoder's hidden states or context embeddings also leads to high computational complexity (Bello et al., 2016).

In this paper, we propose GAA-PtrNet, a novel PtrNet based on graph attention aggregation (GAA) for one-shot DAG scheduling, as shown in Fig 1b. GAA-PtrNet consists of a trainable network to generate pair-wise node attention scores in one shot, and a scheduler that aggregates attention scores to obtain sampling probabilities. In the network, a GNN encoder is employed to obtain DAG task nodes' embeddings, then the pair-wise attention scores between all task nodes are obtained by graph attention mechanism. The attention scores are obtained through a one-shot forward propagation of the network. In the scheduler, at each sequencing step the attention scores between the subgraph formed by the scheduled tasks and each candidate task are aggregated to calculate the probability to select candidate nodes. Furthermore, to train GAA-PtrNet, we design a training strategy based on policy gradient RL with dense reward signal and demonstration learning. **The key contributions of this study are as follows:**



(a) Existing PtrNet.

(b) Our proposed GAA-PtrNet.

Figure 1: The illustrated comparison of existing PtrNet and our GAA-PtrNet in DAG scheduling. $v_{\pi(1)}, v_{\pi(2)}, \dots$ represent the nodes scheduled at each steps.

- 1. We propose GAA-PtrNet, a novel PtrNet based on GAA for one-shot DAG scheduling. By calculating the attention scores among the DAG tasks in a one-shot way, and further computing the decision sampling probability through attention aggregation. GAA-PtrNet has strong capabilities to deal with workflow's topological structures with low time complexity. To the best of our knowledge, this is the first network model to achieve PtrNet-based one-shot DAG scheduling.
- 2. We design a training strategy based on policy gradient RL with dense reward signal and demonstration learning To train GAA-PtrNet for DAG scheduling. Comprehensive experimental results show that the proposed method is superior in terms of objective and runs about 10 times faster when compared to previous PtrNet-based methods, and also performs better than other learning-based DAG scheduling methods.

2 PRELIMINARIES

2.1 DAG TASK MODEL

In a typical DAG scheduling problem G=(V,E,X), the node set $V=\{v_1,v_2,...,v_{|V|}\}$ represents the computation tasks. $X=\{x_1,x_2,...,x_{|V|}\}$ is the attribution of each task node, primarily including the computational workload c_i , the output data size b_i , etc. The directed edges $E\subseteq V\times V$ denotes the precedence constraints among tasks. If $(v_i,v_j)\in E$, v_j cannot start until v_i is completed. For the convenience of calculation, a pseudo entrance task node will be created for the whole DAG scheduling problem, with all the nodes without predecessor to be its successors. Each task v_i is associated with a processing time d_i . Some DAG systems require to consider the data transmission time (z_i) between processors. A task node can be ready and executed only after all its predecessors' execution and transmission is finished.

2.2 DAG SCHEDULING

We consider list scheduling: when existing ready tasks and an idle processor, the ready task with the highest priority is immediately selected and assigned to the processor. Determining the priority list $Solution(G) = [v_{\pi(1)}, v_{\pi(2)}, ..., v_{\pi(|V|)}]$ is the core of DAG scheduling, where $\pi: \{1, 2, ..., |V|\} \rightarrow \{1, 2, ..., |V|\}$. Note that Solution(G) is not necessarily the actual execution order of G. Rather, it is a topological order of G, while each valid Solution(G) does correspond to one valid execution order. We regard each step in constructing this topological order as a decision point, while the selectable nodes at each step are define as the current **candidate nodes**. The ultimate goal is to generate a Solution(G) that optimizes the target, such as makespan.

2.3 Pointer network for DAG scheduling.

The studies that use the foundational PtrNet by Bello et al. (2016) for DAG scheduling (Dong et al. (2021) Zhao et al. (2022) Chen & Wang (2024) Li et al. (2022)) apply a LSTM decoder structure to obtain the current states' embeddings, as shown in equation 1, where $h_1, h_2, ..., h_{|V|}$ are the feature vectors of each task node, $h^{(z)}$ is the encoder's output, and $h_t^{(s)}$ is the decoder's hidden state at t. They then calculate the attention scores (u_i^t) among $h_t^{(s)}$ and the task nodes in the candidate node set $S_C^{(t)}$ by equation 2, where a, W_{ref} and W_q are learnable matrices.

$$h^{(z)} = \text{LSTMEncoder}(h_1, h_2, ..., h_{|V|}), h_t^{(s)} = \text{LSTMdecoder}(h_{t-1}^{(s)}, h^{(z)}) \tag{1}$$

$$u_{i}^{t} = \begin{cases} a^{T} \tanh(W_{ref} h_{i} + W_{q} h_{t}^{(s)}), i \in S_{C}^{(t)} \\ -\infty, i \notin S_{C}^{(t)} \end{cases}$$
 (2)

As for those DAG scheduling studies following the PtrNet of Kool et al. (2018) (e.g., Lee et al. (2020) Lee et al. (2021)), rather than using a LSTM-based structure, they construct a context embedding $h_t^{(c)}$ by concatenating the feature vectors of the task node that are chosen in recent decisions, and the global feature vector of G. They then calculate the attention scores at t (u_i^t) of the candidates in $S_C^{(t)}$ according to $h_t^{(c)}$ by equation 3, where W_Q , W_K are learnable matrices and dim is the embedding's dimension.

$$u_i^t = \begin{cases} \frac{(W_Q h_t^{(c)})(W_K h_i)^T}{\sqrt{dim}}, i \in S_C^{(t)} \\ -\infty, i \notin S_C^{(t)} \end{cases}$$
 (3)

After getting the attention scores (u_i^t) of the nodes in $S_C^{(t)}$, these methods would apply a softmax in order to obtain the probability of selecting each task node among $S_C^{(t)}$, as shown in equation 4, where C_{clip} is the constant coefficient to clip the unnormalized log-probabilities.

$$P(v_i|t) = \frac{\exp(C_{clip} \cdot \tanh(u_i^t))}{\sum_{v_i \in V} \exp(C_{clip} \cdot \tanh(u_j^t)))}$$
(4)

3 Proposed method

3.1 MOTIVATION

Traditional PtrNet-based DAG scheduling methods rely on repeatedly computing the decoder hidden states or context embeddings according to the current observed environment state and the previous decisions at every decision step, and calculating the candidate node attention scores in order to obtain differentiable decision probabilities. Thus, these models tend to rely heavily on recent local decisions, at the expense of capturing the DAG's global structural properties. Furthermore, the repeated computation causes high computational complexity.

We believe that these limitations can be addressed by leveraging the graph attention and its aggregation. Considering that the principle of PtrNet lies in computing the candidate node's attention scores according to the current state (represented by the decoder hidden states or context embeddings), such attention can be obtained directly by aggregating the graph attention scores between the already scheduled nodes and the candidate nodes. Compared with the method in the existing PtrNets, graph attention is more effective in capturing the topological features of DAGs. Moreover, the attention scores between graph nodes can be obtained with one shot forward propagation of the network, which reduces the time complexity.

In this section, we propose GAA-PtrNet, a novel PtrNet based on GAA for one-shot DAG scheduling. GAA-PtrNet consists of a trainable network to generate pair-wise node attention scores, and a scheduler that aggregates attention scores to obtain sampling probabilities. In the network, we utilize a GNN encoder to extract the DAG workflow to obtain node embeddings, then we utilize graph attention mechanism to calculate pair-wise attention scores between all task nodes. The embeddings and attention scores are obtained through a one-shot network forward propagation, which has low time complexity. In the scheduler, at each sequencing step in scheduling, the sampling probabilities of each candidate nodes is obtained through GAA between the scheduled tasks and each candidate tasks. We also design a training strategy based on policy gradient RL with dense reward signal and demonstration learning. The overview of the proposed method is illustrated in Fig 2.

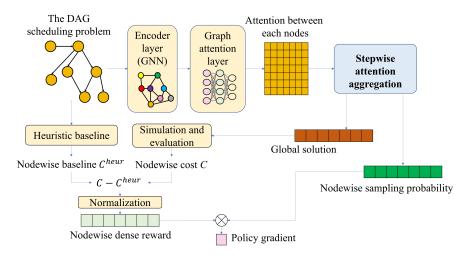


Figure 2: The overall framework of GAA-PtrNet.

3.2 Graph attention score calculation and aggregation

Our proposed GAA-PtrNet consists 2 key components: (1) a network to calculate attention scores in one-shot, and (2) a scheduler to calculate sampling probability through GAA.

We firstly demonstrate the method to obtain pair-wise node raw graph attention scores between task nodes through the network. For each task v_i , its attributions x_i along with environment information are concatenated in to a feature vector, as its raw feature. Then we utilize a bi-directional GNN to extract the DAG's structural information into embedded node features $[h_1, h_2, ..., h_{|V|}]$, as formu-

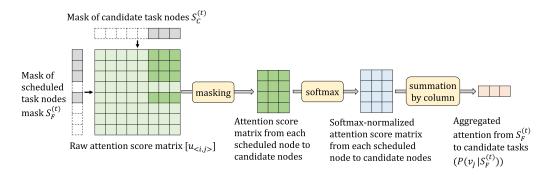


Figure 3: The procedure to obtain sampling probabilities of candidates by GAA. Each term in the raw attention score matrix indicates the attention from a task node (indexed by row) to another (indexed by column). In the attention masks (on the left of the figure), the dark cells indicate the scheduled or candidate nodes. These masks are applied to the raw attention score matrix to obtain the masked attention matrix at t. A softmax operation is then performed over the masked attention (equation 7), followed by a column-wise summation (equation 8), resulting in the aggregated normalized attention from $S_F^{(t)}$ to each in $S_C^{(t)}$, which is the probabilities of selecting each candidates.

lated in equation 5. Subsequently, we compute pair-wise node raw graph attention score from the H and the edges E, yielding an attention matrix U (equation 6), where the attention between nodes v_i and v_j is denoted by $u_{\langle i,j \rangle}$.

$$H = [h_1, h_2, ..., h_{|V|}] = GNN(G)$$
(5)

$$U = [u_{\langle i,j \rangle}] = \operatorname{GraphAttn}(H, E)$$
 (6)

Then, the scheduler computes the decision probability at each step directly according to the graph attention scores. The attention score $\alpha_{< i,j>}^{(t)}$ of attending $S_F^{(t)}$ to each candidate task node $v_j \in S_C^{(t)}$ is calculated by applying a softmax operation over the raw attention scores at the level of the full sets between $S_F^{(t)}$ and $S_C^{(t)}$, as shown in equation 7.

$$\alpha_{\langle i,j\rangle}^{(t)} = \frac{\exp(u_{\langle i,j\rangle})}{\sum_{v_k \in S_C^{(t)}, v_l \in S_F^{(t)}} \exp(u_{\langle k,l\rangle})}, v_i \in S_F^{(t)}, v_j \in S_C^{(t)}$$
(7)

For each candidate node $v_j \in S_C^{(t)}$, we aggregate its attention scores from $S_F^{(t)}$ by summation, as formulated in equation 8. This yields $\alpha_{< S_F^{(t)}, j>}$, the attention from whole subgraph $S_F^{(t)}$ to v_j , which also corresponds to the probability $P(v_j|S_F^{(t)})$ of selecting v_j from $S_C^{(t)}$ at time step t. This probability is differentiable and will be further used in the loss function computation. According to the probability distribution in equation 8, we would sample the node $v_{\pi(t)}$ from $S_C^{(t)}$, as the decision at t. This procedure is illustrated in Fig 3. Note that as established in Section 3, we have assumed that each DAG workflow is augmented with a dummy entrance task node. Hence, at time step 0, when no nodes have been scheduled yet $(S_F^{(0)} = \emptyset)$, attention aggregation is unnecessary. In this occasion, the entry node is selected directly, with $P(v_{\pi(1)}|S_F^{(1)}) = 1$.

$$P(v_j|S_F^{(t)}) = \alpha_{\langle S_F^{(t)}, j \rangle} = \sum_{v_i \in S_F^{(t)}} \alpha_{\langle i, j \rangle}^{(t)}$$
(8)

The softmax result $\alpha_{< i, j>}^{(t)}$ in equation 7 can also be interpreted as the probability to select the node pair < i, j> from all the node pairs between $S_F^{(t)}$ and $S_C^{(t)}$. Then, by performing summation in equation 8, we obtain the joint probability to select all node pairs containing $v_j \in S_C^{(t)}$, which can be interpreted as the probability to select v_j from $S_C^{(t)}$.

equation 7 regards that all raw attention scores are in a unified scale, but the queries for these attention scores are from different task nodes in $S_F^{(t)}$. This raises a potential concern: could differences in the scale of the attention scores across different queries introduce bias in the probabilities computed by equation 8. We believe this issue can be mitigated. In equation 15, the parameters a and W are shared, and the node embeddings H are generated by the same GNN, ensuring that the resulting attention scores lie within a comparable scale. However, since the attention scores are computed in a query-wise manner (i.e., for each node individually), we recommend applying a normalization to the attention scores originating from the same query node $v_i \in S_F^{(t)}$ before applying equation 7. As for equation 14, it adopts a global self-attention, where all attention scores are computed within a unified scale, and thus do not require additional normalization.

3.3 GAA-PTRNET-BASED DAG SCHEDULING PROCESS AND TIME COMPLEXITY ANALYSIS.

The scheduling procedure is demonstrated in Algorithm 1 in Appendix A. Given the scheduling problem G, the network firstly calculates H, then the scheduler updates $S_C^{(t)}$ and $S_F^{(t)}$ step by step, and repeatedly aggregates attention to obtain the sampling probabilities over $S_C^{(t)}$. At each step, a $v_{\pi(t)}$ is sampled accordingly, finally obtaining the whole priority list Solution(G). An example about the procedure for GAA-PtrNet to schedule a DAG is also illustrated in Appendix B.

We analyze the time complexity for GAA-PtrNet to schedule $G=\{V,E,X\}$. Assume the embedding dimension in the network is dim, and the average size of $S_F^{(t)}$ and $S_C^{(t)}$ to be $\overline{|S_F|}$ and $\overline{|S_C|}$. An L-level GNN requires $O(L(|V|\times dim^2+|E|\times dim))$ to generate node embeddings (take graph convolution network for example). Since different PtrNet-based scheduling methods may use similar GNN encoder, we mainly compare the time to make a complete schedule when the embeddings are already obtained. It takes $O(|E|\times dim+|V|\times dim^2)$ for equation 15 or $O(|V|^2+|V|\times|E|+|V|^2\times dim+|V|\times dim^2)$ for equation 14 to compute the raw attention scores. At each timestep, it takes $O(\overline{|S_F|}\times \overline{|S_C|})$ to conduct softmax, and probability calculation is conducted with with $O(\overline{|S_F|}+\overline{|S_C|})$ (upper bounded by O(|V|)). As a brief summary, the time complexity for GAA-PtrNet to make a schedule is $O(|E|\times dim+|V|\times dim^2+|V|\times (\overline{|S_F|}\times \overline{|S_C|}))$ or $O(|V|^2\times dim+|V|\times dim^2+|V|(|V|+\overline{|S_F|}\times \overline{|S_C|}))$, which are both upper bounded by $O(|V|^2\times dim+|V|\times dim^2+|V|^3)$.

Table 1: Time complexity of GAA-PtrNet and existing PtrNet when scheduling DAGs.

Method	Scheduling	Upper bound of scheduling
	$O(E \times dim + V \times dim^2 + V (\overline{ S_F } \times \overline{ S_C }))$	
	$O(V ^2 \times dim^2 + V \times \overline{ S_C } \times dim^2 + V \times \overline{ S_C })$	$O(V ^2 \times dim^2)$
PtrNet-CE	$O(V \times \overline{ S_C }^2 \times dim + V \times \overline{ S_C } \times dim^2)$	$O(V ^2 \times dim^2 + V ^3 \times dim)$

We denote the existing PtrNet base on LSTM structure and additive attention as **PtrNet-LSTM**, and the PtrNet structure with context embedding (CE) and dot-product attention as **PtrNet-CE**. Given the node embeddings, PtrNet-LSTM spends $O(|V|(dim^2 + \overline{|S_C|} \times dim^2))$ to generate a complete solution, in which $O(|V| \times dim^2)$ is the time for the LSTM cell in one step, and $O(\overline{|S_C|} \times dim^2)$ is additive attention. So the overall time complexity for PtrNet-LSTM is $O(|V|^2 \times dim^2 + |V||S_c| \times dim^2 + |V| \times |S_C|)$ upper bounded by $O(|V|^2 \times dim^2)$. Similarly, PtrNet-CE has a $O(|V| \times |S_C|^2 \times dim + |V||S_C| \times dim^2)$ time complexity, which is upper bounded by $O(|V|^2 \times dim^2 + |V|^3 \times dim)$, which is much higher than our proposed GAA-PtrNet. Moreover, in our method, the perstep computational complexity within the iterative scheduling loop is independent of the network dimension dim, as it involves only attention operations. Therefore, the advantage of our approach becomes more pronounced as the number of nodes in the DAG workflow increases. We intuitively compare the scheduling time complexity of GAA-PtrNet and existing PtrNet in Table 1.

3.4 Training strategy with policy gradient RL

We use policy gradient to train our model. Solution(G) and the corresponding differentiable decision probabilities obtained by GAA-PtrNet can form a complete sampling trajectory of RL. Instead

of using the a single final makespan value (C(Solution(G))) of Solution(G) as the reward of this whole sampling trajectory, we introduce dense returns of each individual node-level selection decision as the dense reward signal, in order to guide the optimization of each node-level decision. It is based on the distance of each node's contribution in the objective to the final objective. The return $R_{(t)}$ at time step t can be estimated by equation 9, where $C(v_{\pi(t)})$ indicates the latest completion time among all scheduled nodes when $v_{\pi(t)}$ is finished. The value of $C(v_{\pi(t)})$ can be obtained during the simulation procedure when computing Solution(G).

$$R_{(t)} = C(Solution(G)) - C(v_{\pi(t)})$$
(9)

We further consider introducing advantage baseline on $R_{(t)}$, in order to stabilize the training. The advantage baseline is computed by a heuristic algorithm. Specifically, before training on the workflow instance G, we schedule G by using the heuristic baseline, and record the objective cost value upon the completion of each node $v_{\pi(t)}$, denoted as $C^{heur}(v_{\pi(t)})$. Based on this, we can define the advantage function $A_{(t)}$ as shown in equation 10. $A_{(t)}$ can be further normalized batch-wise into $A^{normalized}_{(t)}$ when there are multiple workflow instances $G = G_1, G_2, ..., G_B$ in training. In this way, the policy gradient loss function can be modified into equation 11, where B is the batch size, and θ represents all the parameters in the trainable neural network.

$$A_{(t)} = R_{(t)} - \mathbf{baseline} = C^{heur}(v_{\pi(t)}) - C(v_{\pi(t)})$$
(10)

$$\nabla_{\theta} J(\theta) = \frac{1}{B} \frac{1}{|V_m|} \sum_{j=1}^{B} \sum_{t=1}^{|V_m|} \nabla_{\theta} \log P(v_{\pi(t)}|S_F^{(t)}) \cdot A_{m,(t)}^{\text{normalized}}$$
(11)

Demonstration learning is applied to initialize the training process, because the trajectories are sampled in a Monte Carlo way: it is generated in one-shot, and no greedy rule is applied in sampling. So, there might be blind exploration at the beginning of the training. We utilize genetic algorithm to generate demonstration solutions to the scheduling problem instances at first, obtaining some suboptimal solutions in the form of task execution orders. These solutions are converted to RL sampling trajectories and the network would be trained on these trajectories for several episodes. This can prevent the network from conducting inefficient exploration at the beginning of training. For detailed description and interpretability analysis of the training strategy, please check Appendix G.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUPS

In this section, we report experimental results to evaluate the contributions of our proposed method, and the performance of our method compared with existing method. The training and simulation of the experiments are conducted on a computer with Ubuntu 20.04 OS, Intel 6226R CPU, 256GB RAM, and RTX 3090 (24GB) graphic card. The experiments are conducted using Python 3.9 as the programming language. The neural network model is implemented based on PyTorch 2.7 and torch-geometric 2.6. Details about the implementation of experiments are presented in Appendix E.

In order to thoroughly evaluate the adaptability and robustness of our proposed method across different scenarios, we introduce the following benchmarks in our experiments: **TPC-H** represents the real-world DAG workflow scheduling scenario under homogeneous processor environment. Here, we use the TPC-H benchmark generated by Wang et al. (2021). **Pegasus** is a real-world scientific workflow scheduling tracing dataset with heterogeneous multiprocessor environment (Deelman et al., 2015). **Randomly generated DAG workflows**, **TPC-H** and **Pegasus**. In the **randomly generated DAG workflows**, we generate DAG scheduling problems with varying shapes, sizes, and task node attributes by tuning the parameters, following the DAG generation paradigm of Topcuoglu et al. (2002). These parameters include graph shape parameter (β_1), task node heterogeneity (β_2), Computation-to-Communication Ratio (CCR) and average number of tasks in each sub-DAG. Please check Appendix D for more details about these benchmarks.

The target mainly includes average **makespan** of the DAG workflows. Also, **speedup** and **relative gap** are applied to evaluate the performance. Speedup indicates how many times the generated solution is faster in makespan to the case when all tasks are executed only on the fastest processor. Relative gap indicates the gap in makespan relative to the heuristic baseline. Besides, we test the **runtime** to infer a solution, in order to evaluate the time complexity of our method in practice.

4.2 RESULTS AND DISCUSSION

4.2.1 ABLATION STUDY ABOUT GAA-PTRNET

We compare the performance of GAA-PtrNet on DAG scheduling to **PtrNet-LSTM** and **PtrNet-CE**. We also evaluate the influence of different attention mechanism used in GAA-PtrNet: the GAA-PtrNet implemented by self-product attention with position encoding (Ying et al., 2021), donated as **GAA-PtrNet-SA**, and the implementation based on classic graph attention network (Brody et al., 2021), donated as **GAA-PtrNet-GAT**. Note that these are originally for graph representation propose, not for scheduling, modifications are necessary, please check Appendix E.1 for more details.

As presented in Table 3 and Table 2, and the additional results in Appendix F, both **GAA-PtrNet-SA** and **GAA-PtrNet-GAT** outperform **PtrNet-CE** and **PtrNet-LSTM** in different evaluation matrices on diverse benchmarks. In most cases, the performance difference between **GAA-PtrNet-SA** and **GAA-PtrNet-GAT** is minor. This shows that the key factor behind the performance improvement is our proposed GAA, rather than which graph attention computation method is selected.

We evaluated the average runtime of the model, i.e., the time required to infer a solution for a DAG scheduling instance (Table 4). As predicted by the time complexity analysis, the one-shot scheduling by GAA-PtrNet exhibits substantially better runtime performance than PtrNet-LSTM and PtrNet-CE. Our method averagely runs 10 time faster. We further observed that problem size |V| is the dominant factor influencing runtime, while the choice of graph attention computation method has minor effect. We present only representative results on Pegasus in Table 4, since we found that the runtime difference on different benchmarks is also minor. This shows that the runtime of GAA-PrtNet is less sensitive to the DAG topology: although the number of edges |E| also affects the complexity in the theoretical analysis, the main factor remains |V|. We believe this is because GAA-PrtNet must explicitly invoke $\left|V\right|$ times of loops to produce a complete solution, which is the most time-consuming operation. Moreover, it can be observed that the actual runtime of our GAA-PtrNet-based one-shot scheduling method increases relatively slowly with the growth of the scheduling problem size, unlike existing PtrNet models whose runtime scales linearly with problem size. We attribute this to the fact that conventional PtrNets repeatedly compute hidden states (or context embeddings) and attention at every scheduling decision, which accounts for the majority of the runtime, whereas our method avoids this overhead.

4.2.2 Comparison with baseline methods

Our approach is compared against these baselines: (1) The heuristic algorithms used as advantage baselines in our method; (2) Jeon et al. (2023), a state-of-the-art RL-based one-shot DAG scheduling method; (3) **EGS** (Sun et al., 2024), a DAG scheduling approach based on edge generation; (4) **POMO-DAG**, our adapted implementation of the POMO (Kwon et al., 2020) for DAG scheduling. Check Appendix E.3 for more details about the implementation of these baseline methods.

The results of comparison experiments are presented in Table 3, Table 2 and Appendix F. Our method outperforms the RL-based one-shot DAG scheduling method proposed by Jeon et al. (2023). We argue that this is because their method relies on the Plackett–Luce ranking model to generate a priority list of DAG nodes, which is highly sensitive to the numerical distribution of the assigned logits (i.e., priority values). This results in less stability in training and may require extensive hyperparameter tuning. On the contrary, our method does not need to conduct such ranking model. Since Jeon et al. (2023) is a one-shot scheduling algorithm, we also compared the runtime of their method with ours in Table 4. The results show that our method is not slower; on the contrary, it is even faster on problems of scale 400. When compared with **POMO-DAG**, our method consistently yields better results. Compared to **EGS**, our method achieves similar or even better solution quality. We attribute this to that the search space in EGS is significantly larger. In contrast, our search space

is constrained to the candidate node set. In summary, our method demonstrates advantages over existing learning-based DAG scheduling methods across different DAG scheduling scenarios.

Table 2: Experimental results on TPC-H benchmark, with 3 different problem instance sizes (50, 100 and 150 sub-DAGs, with each sub-DAG containing averagely 9.17 task nodes)

Method	1	ГРС-Н 50)	T	PC-H 10	0	TPC-H 150			
	makespan gap/% speed		speed up	makespan	gap/%	speed up	makespan	gap/%	speed up	
GAA-PtrNet-SA GAA-PtrNet-GAT	21.37 21.33	-14.42 -14.58	5.25 5.26	39.59 42.18	-7.54 -1.49	5.37 5.04	67.08 67.81	-3.84 -2.79	4.81 4.96	
PtrNet-LSTM PtrNet-CE	26.10 26.19	4.53 4.88	4.30 4.28	44.94 44.18	4.95 3.18	4.73 4.81	73.20 72.84	4.93 4.42	4.41 4.43	
STF (heuristic baseline) Jeon et al. (2023) POMO-DAG EGS	24.97 23.73 46.90 24.58	-4.95 87.8 -1.55	4.49 4.73 2.39 4.56	42.85 41.22 90.30 42.18	-3.81 110.74 -1.56	4.96 5.15 2.35 5.04	69.76 74.02 141.90 68.99	6.11 103.43 -1.10	4.63 4.35 2.27 4.68	

Table 3: Experimental results on SIPHT dataset in Pegasus, with 4 different problem instance sizes (averagely 100, 200, 300, and 400 task nodes).

Method	SIPHT-100			SIPHT-200			SI	PHT-30)0	SIPHT-400			
	makespan	gap/%	speed up	makespan	gap/%	speed up	makespan	gap/%	speed up	makespan	gap/%	speed up	
GAA-PtrNet-SA GAA-PtrNet-GAT	191.1 191.7	-15.81 -15.55	2.43 2.42	340.3 338.8	-4.891 -3.41	2.45 2.51	542.1 542.7	-0.20 -0.05	2.51 2.50	708.5 708.3	-0.88 -0.91	2.51 2.51	
PtrNet-LSTM PtrNet-CE	205.0 207.5	-9.69 -8.59	2.27 2.24	352.8 354.8	-1.40 -0.84	2.41 2.40	552.2 557.0	1.69 2.58	2.46 2.44	717.5 719.6	0.38 0.67	2.48 2.476	
HEFT (heuristic baseline) Jeon et al. (2023) POMO-DAG EGS	227.0 218.5 214.0 200.6	-3.74 -5.74 -11.63	2.05 2.23 2.17 2.32	357.8 352.2 367.5 346.3	-1.57 2.72 -3.21	2.38 2.42 2.31 2.46	543.0 550.6 575.8 542.8	1.40 6.05 -0.04	2.51 2.47 2.36 2.51	714.8 712.7 741.9 710.2	-0.29 3.80 -0.42	2.49 2.50 2.40 2.51	

Table 4: Runtime comparison on Pegasus. We tested DAGs of various scales and, for each scale, compared results across different benchmarks. Each value in the table denotes the runtime (in seconds).

Method		size =	100		size = 1	200		size = 1	300	size = 400			
	SIPHT	LIGO	GENOME	SIPHT	LIGO	GENOME	SIPHT	LIGO	GENOME	SIPHT	LIGO	GENOME	
GAA-PtrNet-SA GAA-PtrNet-GAT	0.126		0.130 0.129	0.219 0.228	0.226 0.226	0.225 0.223		0.300 0.297	0.299 0.298	0.347 0.355	0.356 0.357	0.354 0.356	
PtrNet-LSTM PtrNet-CE	1.343	1.389 1.097	1.387 1.097	2.645 2.081		2.710 2.130	3.921 3.086	4.054 3.188	4.015 3.159	5.202 4.108		3.994 3.156	
Jeon et al. (2023)	0.07	0.08	0.08	0.15	0.16	0.15	0.26	0.26	0.27	0.43	0.43	0.46	

5 Conclusion

In this paper, we propose a novel PtrNet based on GAA for one-shot DAG scheduling. By calculating the graph attention scores in a one-shot way and obtaining the task node sampling probability through GAA, GAA-PtrNet achieves one-shot DAG scheduling. It can handle with DAG workflows' complex topological structure with low time complexity. We also propose a RL-based policy gradient training strategy for GAA-PtrNet. We conducted comparative and ablation experiments on various DAG scheduling scenarios, demonstrating the superiority of our method in terms of solution quality and runtime. This suggest that our method have potential for further extension and optimization in large-scale or real-time DAG workflow environments. This is the first network model to achieve PtrNet-based one-shot DAG scheduling. As future work, we will expand this foundational study to more specific high performance computation applications by considering more domain knowledge and multiple optimization objectives.

6 ETHICS AND REPRODUCIBILITY STATEMENT

We have read ICLR code of ethics and we claim no potential violations of the code of ethics. Regarding reproducibility, we describe the experimental environment in Section 4, provide details of the benchmark sources in Appendix D, and present the implementation details in Appendix E.

REFERENCES

- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- Irwan Bello, Hieu Pham, Quoc V Le, Mohammad Norouzi, and Samy Bengio. Neural combinatorial optimization with reinforcement learning. *arXiv preprint arXiv:1611.09940*, 2016.
- Shaked Brody, Uri Alon, and Eran Yahav. How attentive are graph attention networks? *arXiv* preprint arXiv:2105.14491, 2021.
- Genxin Chen, Jin Qi, Ying Sun, Xiaoxuan Hu, Zhenjiang Dong, and Yanfei Sun. A collaborative scheduling method for cloud computing heterogeneous workflows based on deep reinforcement learning. *Future Generation Computer Systems*, 141:284–297, 2023. ISSN 0167-739X. doi: 10.1016/j.future.2022.11.032. URL https://www.sciencedirect.com/science/article/pii/S0167739X22004034.
- Suhong Chen and Xudong Wang. Ptrtasking: Pointer network based task scheduling for multi-connectivity enabled mec services. *IEEE Transactions on Mobile Computing*, 23(10):9398–9409, 2024. doi: 10.1109/TMC.2024.3363662.
- Ewa Deelman, Karan Vahi, Gideon Juve, Mats Rynge, Scott Callaghan, Philip J. Maechling, Rajiv Mayani, Weiwei Chen, Rafael Ferreira da Silva, Miron Livny, and Kent Wenger. Pegasus, a workflow management system for science automation. *Future Generation Computer Systems*, 46: 17–35, 2015. ISSN 0167-739X. doi: https://doi.org/10.1016/j.future.2014.10.008. URL https://www.sciencedirect.com/science/article/pii/S0167739X14002015.
- Michel Deudon, Pierre Cournut, Alexandre Lacoste, Yossiri Adulyasak, and Louis-Martin Rousseau. Learning heuristics for the tsp by policy gradient. In *Integration of Constraint Programming, Artificial Intelligence, and Operations Research: 15th International Conference, CPAIOR 2018, Delft, The Netherlands, June 26–29, 2018, Proceedings 15*, pp. 170–181. Springer, 2018.
- Hamza Djigal, Jun Feng, Jiamin Lu, and Jidong Ge. Ippts: An efficient algorithm for scientific workflow scheduling in heterogeneous computing systems. *IEEE Transactions on Parallel and Distributed Systems*, 32(5):1057–1071, 2021. doi: 10.1109/TPDS.2020.3041829.
- Tingting Dong, Fei Xue, Changbai Xiao, and Jiangjiang Zhang. Deep reinforcement learning for dynamic workflow scheduling in cloud environment. In 2021 IEEE International Conference on Services Computing (SCC), pp. 107–115, 2021. doi: 10.1109/SCC53864.2021.00023.
- Tingting Dong, Fei Xue, Hengliang Tang, and Chuangbai Xiao. Deep reinforcement learning for fault-tolerant workflow scheduling in cloud environment. *Applied Intelligence*, 53(9): 9916-9932, 2023. ISSN 0924-669X. doi: 10.1007/s10489-022-03963-w. URL https://link.springer.com/article/10.1007/s10489-022-03963-W.
- Yan Gu, Zhaoze Liu, Shuhong Dai, Cong Liu, Ying Wang, Shen Wang, Georgios Theodoropoulos, and Long Cheng. Deep reinforcement learning for job scheduling and resource management in cloud computing: An algorithm-level review. *arXiv preprint arXiv:2501.01007*, 2025.
- Mirsaeid Hosseini Shirvani. A survey study on task scheduling schemes for workflow executions in cloud computing environment: classification and challenges. *The Journal of Supercomputing*, 80 (7):9384–9437, 2024.
- Wonseok Jeon, Mukul Gagrani, Burak Bartan, Weiliang Will Zeng, Harris Teague, Piero Zappi, and Christopher Lott. Neural dag scheduling via one-shot priority sampling. In *The Eleventh International Conference on Learning Representations*, 2023.

- Athanassios M. Kintsakis, Fotis E. Psomopoulos, and Pericles A. Mitkas. Reinforcement learning based scheduling in a workflow management system. *Engineering Applications of Artificial Intelligence*, 81:94–106, 2019. ISSN 0952-1976. doi: https://doi.org/10.1016/j.engappai. 2019.02.013. URL https://www.sciencedirect.com/science/article/pii/S0952197619300351.
 - Wouter Kool, Herke Van Hoof, and Max Welling. Attention, learn to solve routing problems! *arXiv* preprint arXiv:1803.08475, 2018.
 - Yeong-Dae Kwon, Jinho Choo, Byoungjip Kim, Iljoo Yoon, Youngjune Gwon, and Seungjai Min. Pomo: Policy optimization with multiple optima for reinforcement learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neu*ral Information Processing Systems, volume 33, pp. 21188–21198. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/ file/f231f2107df69eab0a3862d50018a9b2-Paper.pdf.
 - Hyunsung Lee, Jinkyu Lee, Ikjun Yeom, and Honguk Woo. Panda: Reinforcement learning-based priority assignment for multi-processor real-time scheduling. *IEEE Access*, 8:185570–185583, 2020. doi: 10.1109/ACCESS.2020.3029040.
 - Hyunsung Lee, Sangwoo Cho, Yeongjae Jang, Jinkyu Lee, and Honguk Woo. A global dag task scheduler using deep reinforcement learning and graph convolution network. *IEEE Access*, 9: 158548–158561, 2021. doi: 10.1109/ACCESS.2021.3130407.
 - Huifang Li, Jianghang Huang, Binyang Wang, and Yushun Fan. Weighted double deep q-network based reinforcement learning for bi-objective multi-workflow scheduling in the cloud. *Cluster Computing*, 25(2):751–768, 2022.
 - Qiang Ma, Suwen Ge, Danyang He, Darshan Thaker, and Iddo Drori. Combinatorial optimization by graph pointer networks and hierarchical reinforcement learning. *arXiv* preprint arXiv:1911.04936, 2019.
 - Hongzi Mao, Malte Schwarzkopf, Shaileshh Bojja Venkatakrishnan, Zili Meng, and Mohammad Alizadeh. Learning scheduling algorithms for data processing clusters. In *Proceedings of the ACM Special Interest Group on Data Communication*, SIGCOMM '19, pp. 270–288, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450359566. doi: 10. 1145/3341302.3342080. URL https://dl.acm.org/doi/abs/10.1145/3341302.3342080.
 - Alexander Nasuta, Marco Kemmerling, Daniel Lütticke, and Robert H Schmitt. Reward shaping for job shop scheduling. In *International Conference on Machine Learning, Optimization, and Data Science*, pp. 197–211. Springer, 2024. ISBN 978-3-031-53969-5.
 - Xumai Qi, Dongdong Zhang, Taotao Liu, and Hongcheng Wang. Deep reinforcement learning for large-scale scientific workflow scheduling with improved structure feature extraction and sampling. In *IFIP International Conference on Network and Parallel Computing*, pp. 310–323. Springer, 2024.
 - Shuo Qin, Dechang Pi, Zhongshi Shao, Yue Xu, and Yang Chen. Reliability-aware multi-objective memetic algorithm for workflow scheduling problem in multi-cloud system. *IEEE Transactions on Parallel and Distributed Systems*, 34(4):1343–1361, 2023.
 - Wen Shi and Chengpu Yu. Multi-agent task allocation with multiple depots using graph attention pointer network. *Electronics*, 12(16), 2023. ISSN 2079-9292. URL https://www.mdpi.com/2079-9292/12/16/3378.
 - Yutao Song, Chen Li, Lihua Tian, and Hui Song. A reinforcement learning based job scheduling algorithm for heterogeneous computing environment. *Computers and Electrical Engineering*, 107:108653, 2023. ISSN 0045-7906. doi: 10.1016/j.compeleceng.2023.108653. URL https://www.sciencedirect.com/science/article/pii/S0045790623000782.
 - Binqi Sun, Mirco Theile, Ziyuan Qin, Daniele Bernardini, Debayan Roy, Andrea Bastoni, and Marco Caccamo. Edge generation scheduling for dag tasks using deep reinforcement learning. *IEEE Transactions on Computers*, 73(4):1034–1047, 2024. doi: 10.1109/TC.2024.3350243.

- H. Topcuoglu, S. Hariri, and Min-You Wu. Performance-effective and low-complexity task scheduling for heterogeneous computing. *IEEE Transactions on Parallel and Distributed Systems*, 13(3): 260–274, 2002. doi: 10.1109/71.993206.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 28. Curran Associates, Inc., 2015. URL https://proceedings.neurips.cc/paper_files/paper/2015/file/29921001f2f04bd3baee84a12e98098f-Paper.pdf.
- Runzhong Wang, Zhigang Hua, Gan Liu, Jiayi Zhang, Junchi Yan, Feng Qi, Shuang Yang, Jun Zhou, and Xiaokang Yang. A bi-level framework for learning to solve combinatorial optimization on graphs. 34:21453-21466, 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/b2f627ffff19fda463cb386442eac2b3d-Paper.pdf.
- Xiao Wang, Hanchuan Xu, Xianzhi Wang, Xiaofei Xu, and Zhongjie Wang. A graph neural network and pointer network-based approach for qos-aware service composition. *IEEE Transactions on Services Computing*, 16(3):1589–1603, 2023. doi: 10.1109/TSC.2022.3196915.
- Xiaohan Wang, Yuanjun Laili, Lin Zhang, and Yongkui Liu. Hybrid task scheduling in cloud manufacturing with sparse-reward deep reinforcement learning. *IEEE Transactions on Automation Science and Engineering*, 22:1878–1892, 2025. doi: 10.1109/TASE.2024.3371250.
- Yi Xie, Feng-Xian Gui, Wei-Jun Wang, and Chen-Fu Chien. A two-stage multi-population genetic algorithm with heuristics for workflow scheduling in heterogeneous distributed computing environments. *IEEE Transactions on Cloud Computing*, 11(2):1446–1460, 2021.
- Zhe Yang, Phuong Nguyen, Haiming Jin, and Klara Nahrstedt. Miras: Model-based reinforcement learning for microservice resource allocation over scientific workflows. In 2019 IEEE 39th International Conference on Distributed Computing Systems (ICDCS), pp. 122–132, 2019. doi: 10.1109/ICDCS.2019.00021.
- Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-Yan Liu. Do transformers really perform badly for graph representation? In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), Advances in Neural Information Processing Systems, volume 34, pp. 28877–28888. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/f1c1592588411002af340cbaedd6fc33-Paper.pdf.
- Xiaoming Yu, Wenjun Wu, and Yangzhou Wang. Integrating cognition cost with reliability qos for dynamic workflow scheduling using reinforcement learning. *IEEE Transactions on Services Computing*, 16(4):2713–2726, 2023. doi: 10.1109/TSC.2023.3253182.
- Yuqi Zhao, Bing Li, Jian Wang, Delun Jiang, and Duantengchuan Li. Integrating deep reinforcement learning with pointer networks for service request scheduling in edge computing. *Knowledge-Based Systems*, 258:109983, 2022. ISSN 0950-7051. doi: https://doi.org/10.1016/j.knosys. 2022.109983. URL https://www.sciencedirect.com/science/article/pii/S0950705122010760.
- Yunfan Zhou, Xijun Li, Jinhong Luo, Mingxuan Yuan, Jia Zeng, and Jianguo Yao. Learning to optimize dag scheduling in heterogeneous environment. In 2022 23rd IEEE International Conference on Mobile Data Management (MDM), pp. 137–146, 2022. doi: 10.1109/MDM55031. 2022.00040.
- Zhaomeng Zhu, Gongxuan Zhang, Miqing Li, and Xiaohui Liu. Evolutionary multi-objective workflow scheduling in cloud. *IEEE Transactions on Parallel and Distributed Systems*, 27(5):1344–1357, 2016. doi: 10.1109/TPDS.2015.2446459.

THE ALGORITHM OF ONE-SHOT DAG SCHEDULING BY GAA-PTRNET

Algorithm 1 Scheduling a DAG workflow instance by GAA-PtrNet

Input: Workflow instance G = (V, E, I)

Output: Task node priority list $[v_{\pi(1)}, v_{\pi(2)}, ..., v_{\pi(|V|)}]$

- 1: Calculate the node embeddings H by equation 5
- 2: Calculate the raw attention scores $\{u_{\langle i,j\rangle}\}$ by equation 15 or equation 14 on $G_{reserved}$.
- 3: Set the time step t=1, the priority list $O=[\]$, the scheduled set $S_E^{(t)}=\emptyset$
- 4: while $t \leq |V|$ do
- Collet the current candidate task node set $S_C^{(t)}$ 5:
- Softmax the raw attention scores over $S_C^{(t)}$ and $S_F^{(t)}$ by equation 7 Calculate the probability to select each node in $S_C^{(t)}$ by equation 8 6:
- 7:
- Sample $v_{\pi(t)}$ based on the probabilities. Append $v_{\pi(t)}$ to O and insert it to $S_F^{(t)}$ 8:
- 9: Assign $v_{\pi(t)}$ to a processor according to the EFT principle.
- 10: end while

A ILLUSTRATIVE EXAMPLE OF ONE-SHOT DAG SCHEDULING BY **GAA-PTRNET**

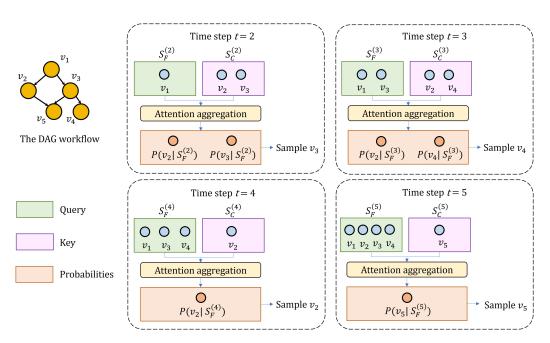


Figure 4: An example of scheduling a DAG with 5 task nodes by GAA-PtrNet

Fig 4 shows an example of scheduling a DAG with 5 task nodes by GAA-PtrNet. In Fig 4, at the beginning, the pseudo entry node v_1 is processed in advance. At timestep 2, the scheduled node subset is $S_F^{(2)} = v_1$, and the candidate subset is $S_C^{(2)} = v_2, v_3$. Through GAA, we obtain the probability to sample v_2 and v_3 , and the node that the scheduler eventually sample is v_3 . Next, at timestep 3, the scheduled node subset is $S_F^{(3)} = v_1, v_3$, and the candidate subset is $S_C^{(2)} = v_2, v_4$. Assume v_4 is sampled. Then, similarly, the scheduler samples v_2 at timestep 4, and v_5 at timestep 5. The final priority list is $Solution(G) = [v_1, v_3, v_4, v_2, v_5]$.

C RELATED WORKS

702

703 704

705706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727 728

729 730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747 748 749

750 751

752 753

754

755

C.1 RL FOR DAG SCHEDULING

While most existing studies on DAG workflow scheduling in cloud computing and cluster computing area still rely on heuristic (Topcuoglu et al., 2002) (Djigal et al., 2021) or meta-heuristic algorithms (Xie et al., 2021) (Qin et al., 2023), learning-based approaches are gradually becoming the mainstream. Yang et al. (2019) was the first to introduce model-based reinforcement learning in the distributed and cloud-based system to schedule scientific workflow. Mao et al. (2019) first proposed to introduce GNN into RL model in order to extract structural features of the workflow instances to make better decisions. But it cannot handle heterogeneous computing resources and large scale workflows. A GNN encoder to extract workflows' structural information, and a policy network to output scheduling decisions, have been the common network architecture in the RL-based workflow scheduling studies(Zhou et al., 2022) (Song et al., 2023) (Qi et al., 2024). Most of the existing studies followed the Markov Decision Process (MDP) to select a candidate task node at each step, constructing the whole solution incrementally (Yu et al., 2023) (Dong et al., 2023). Jeon et al. (2023) argued that such methods require to re-encode the entire graph instance in each step, leading to excessive computation, so they attempted to learn to output the node logits as their global fixed priorities for all nodes in a one-shot manner. But its training stability was highly sensitive to the distribution of the logits. Wang et al. (2021) and Sun et al. (2024) attempted to learn to prioritize task nodes indirectly by modifying the DAG's topological structure, rather than directly generating task priority lists. But these approaches often suffered from an excessively large search space. A common challenge for RL-based workflow scheduling approaches is sparse reward in training: the agent can obtain an accurate feed back on its decision only after the whole workflow instance is scheduled, so that the model has to be optimized with a limited amount of global reward signal. Some works Chen et al. (2023) Wang et al. (2025) Nasuta et al. (2024) tried to overcome this challenge, but they were limited in specific application domains, rather than providing general solutions.

C.2 RL COMBINED WITH PTRNET FOR DAG SCHEDULING

Just like other COP, workflow scheduling in high performance computing environment involves reordering the components (i.e., task, resource, or meta-heuristics rules) of the given problem instance. Under this background, PtrNet, as a sequence learning method, have shown its distinct advantages. PtrNet was initially proposed to handle sequence-to-sequence learning tasks in natural language processing area (Vinyals et al., 2015). Bello et al. (2016) firstly proposed to utilize Ptr-Net to solve routing problems, and Ma et al. (2019) introduced GNN to PtrNet as the encoder to better encode the graph-structured problem instance. Their basic idea was followed by many DAG workflow scheduling studies Dong et al. (2021) Zhao et al. (2022) Chen & Wang (2024) Li et al. (2022). To address the limitations of LSTM-based PtrNet in handling non-sequential structures, Deudon et al. (2018) modified the method of Bello et al. (2016) by computing the attention of context embedding, instead of the LSTM decoder's hidden state, to the candidate actions. In order to improve the computation parallelization and scaling potential, Kool et al. (2018) further modified the attention calculation in in PtrNet-based COP solver with dot-product attention (Vaswani et al., 2017), instead of the traditional additive attention mechanism (Bahdanau et al., 2014). Their method are followed by many recent task scheduling studies Lee et al. (2020) Lee et al. (2021) (Shi & Yu, 2023) (Wang et al., 2023). But in general, these methods suffer from high computational complexity and are often limited to capturing local information within the graph. Moreover, the potential of GNN-based graph attention mechanisms (Veličković et al., 2017) (Ying et al., 2021) (Brody et al., 2021) for pointer networks has not been fully explored.

D MORE INFORMATION ABOUT THE BENCHMARKS

D.1 RANDOMLY GENERATED DAG WORKFLOWS

To evaluate the performance of our method in scheduling multiple workflow instances, we generate a set of random DAG workflows as test cases. We adopt it as the benchmark for testing DAG workflow scheduling algorithm under the simulated heterogeneous multiprocessor setting. Under such setting, the attributes of task node v_i include the computational workload c_i and the output data size b_i . For

 each processor m, the key attribute is its computational capacity f_m . Assuming that task node v_i is assigned to processor m, the processing time d_i (as described in section 2) can be calculated as Equation equation 12.

$$d_i = \frac{c_i}{f_m} \tag{12}$$

Besides, heterogeneous multiprocessor environment requires to consider the transmission time z_i between processors: a task node cannot begin execution until all of its predecessor nodes have completed both their computation, and the transmission of their output data to the processor on which it is scheduled. Specifically, if a task node v_i and its successor v_j are assigned to different processors, then a transmission time is related to the output data size of v_i and the bandwidth of computing environment. If both v_i and v_j are assigned to the same processor, the transmission time is considered negligible. This definition of z_i is formalized in Equation equation 13.

$$z_i = \begin{cases} \frac{b_i}{\text{bandwidth}}, & \text{if } v_i \text{ and } v_j \text{ are on different processors} \\ 0, & \text{if } v_i \text{ and } v_j \text{ are on the same processor} \end{cases}$$
 (13)

Following Topcuoglu et al. (2002), the randomly generation process considers the following key parameters: (1) **Graph shape parameter** (β_1): this parameter characterizes the depth of the DAG. (2) **Task node heterogeneity** (β_2): this captures the diversity in computational workloads c_i and data sizes b_i among different task nodes. (3) **Computation-to-Communication Ratio** (CCR): defined as the ratio between the average task processing time and the average data transmission time. (4) **Average number of tasks** in each sub-DAG in the whole DAG scheduling problem.

By tuning these parameters, we construct a diverse set of DAG workflows with varying structures and heterogeneity levels, simulating real-world heterogeneous environments in the parallel computing area. In the experimental setup, we vary one parameter at a time while fixing the remaining parameters to randomly assigned constant values. For each configuration, we generate multiple scheduling problem instances to examine how the method's performance changes with respect to the chosen parameter. Specifically, the parameter β_1 is tested with a range of $\beta_1 = \{0.1, 1.0, 1.5\}$. For β_2 , it's tested with a range of $\beta_2 = \{1.0, 1.5\}$. CCR is in a range of $CCR = \{0.1, 0.5, 1.0, 2.0, 5.0, 10.0\}$. |V| is tested in a range of $|V| = \{10, 20, 30\}$. Each scheduling problem instance contains 20 independent sub-DAGs.

Our proposed method outputs only the execution order of the nodes, while the assignment of each node to a processor is determined using the Earliest Finish Time (EFT)-greedy rule. Specifically, for a given node to be scheduled, which is determined by the RL network, the EFT-greedy rule dispatches it to the processor that results in the earliest possible finish time. We adopt HEFT (Topcuoglu et al., 2002) as the advantage baseline for dense reward signal under this setting. HEFT is a classic list-based heuristic algorithm for DAG scheduling on heterogeneous processors. It computes the priority (rank-up) of each node based on the average finish time of its successor nodes, and then assigns each node to a processor using the EFT-greedy rule.

D.2 TPC-H

We adopt TPC-H as the benchmark for DAG workflow scheduling under the homogeneous single-processor setting. Since it's a homogeneous processor scenario, there is no need to dispatch tasks to specific processors in such setting. Each task node v_i has a fixed processing time d_i and computation resource requirement q_i . The total resource consumption of concurrently running tasks must not exceed the system's maximum resource capacity. We use the open source code 1 implemented by Wang et al. (2021) to generate TPC-H instances. Shortest Time First (STF) heuristic is adopted as the advantage baseline for TPC-H in our research. At each decision point, the STF rule selects the task with the shortest processing time.

https://github.com/Thinklab-SJTU/PPO-BiHyb/tree/main/dag_data/tpch

D.3 PEGASUS

Pegasus (Deelman et al., 2015) ² provides an open-source workflow trace data generated from various scientific computing applications. We adopt LIGO, SIPHT and Genome data in Pegasus dataset as the benchmark for DAG workflow scheduling under the real-world heterogeneous multiprocessor setting. Since it's also a heterogeneous multiprocessor environment, similar with the randomly generated workflows, the task processing time is calculated by equation 12 and transmission time should be considered calculated by equation 13. EFT-greedy rule is utilized to dispatch the scheduled task nodes to the processors. HEFT is used as the advantage baseline for the dense reward signal.

E IMPLEMENTATION DETAILS

E.1 IMPLEMENTATION ABOUT THE GRAPH ATTENTION FOR SCHEDULING

We compare GAA-PtrNet with 2 different graph attention calculation methods: the GAA-PtrNet implemented by self-product attention with position encoding in equation 14 (Brody et al., 2021), and the implementation by classic graph attention in equation 15 (Ying et al., 2021), where a, W, W_Q and W_K are learnable weight matrices, and dim is the embedding dimensionality. equation 14 applies self-attention over all node pairs at the graph level, and incorporates a learnable bias term $b_{\phi(v_i,v_j)}$ to encode the shortest-path distance $\phi(v_i,v_j)$ between nodes. equation 15 computes attention only between adjacent nodes, and we mask out non-adjacent pairs.

$$u_{\langle i,j\rangle} = \left[\frac{(HW_Q)(HW_K)^T}{\sqrt{dim}}\right]_{\langle i,j\rangle} + b_{\phi(v_i,v_j)}$$
(14)

$$u_{\langle i,j\rangle} = \begin{cases} a^T W[h_i||h_j] & , (v_j, v_i) \in E \\ -\infty & , (v_j, v_i) \notin E \end{cases}$$

$$(15)$$

It is important to note that in the original formulation of Brody et al. (2021), a LeakyReLU layer is applied to $u_{\langle i,j\rangle}$, because it is then used to compute weights for the following feature aggregation. However, in our method, since we directly use the raw attention scores to compute softmax probabilities, we omit the activation layer in equation 15. This modification ensures that the softmax is directly conducted on the raw attention logits, and the results are more interpretable as probability.

Moreover, in the original paper of Brody et al. (2021) and Ying et al. (2021), the attention follows the direction of edges in the graph. But in GAA-PtrNet, the attention scores is employed to compute the sampling probabilities. Therefore, to compute the attention from the scheduled node set to the unscheduled nodes, the edge directions must be reversed. As a result, we apply these attention module to the reversed DAG, $G_{reserved}$.

E.2 NEURAL NETWORK AND HYPER PARAMETER SETTINGS

In the trainable neural network of GAA-PtrNet, as well as the implementation of **PtrNet-LSTM** and **PtrNet-CE** in ablation study, we use a Graphormer(Ying et al., 2021) GNN 3 (which is released under the MIT license) with 4 layers and 4 attention heads, obtaining node embeddings $[h_1,h_2,...,h_{|V|}]$. It processes both the input DAG and its reversed graph simultaneously and outputs 32-dimension node embeddings for each task nodes. In our proposed method and the ablation study, each attention mechanism used to output the raw attention scores for GAA is restricted to a single head, and the dimension of its learnable matrices is 32, in order to keep consistent with the dimension of the node embeddings. We train each benchmark for at most 5000 epochs, although it actually converged much earlier than this epoch. The batch size is set to be 16. The Adam optimizer is employed with learning rate 5×10^{-4} .

²https://pegasus.isi.edu/workflow_gallery/

³https://github.com/microsoft/Graphormer

E.3 BASELINE IMPLEMENTATION

Jeon et al. (2023). The original authors didn't release their source code, the idea described in their paper is sufficiently clear and straightforward for us to reproduce.

POMO-DAG. We build POMO-DAG upon the POMO ⁴ proposed by Kwon et al. (2020), adapting its problem instance encoder to a Graphormer-based GNN (Ying et al., 2021) so that it can process DAG scheduling problems.

EGS. For EGS (Sun et al., 2024), the basic framework of the original code is publicly available⁵. We retained the original structure and implemented the missing policy network and training procedure that were not released.

E.4 SIMULATION AND EVALUATION ENVIRONMENT

To evaluate the generated scheduling solutions and obtain both the overall optimization objective and node-level dense reward signals, we implemented a DAG workflow simulation environment based on the open-source SimPy ⁶ platform in Python language. This simulator is further wrapped into an OpenAI Gym⁷ environment to integrate with reinforcement learning frameworks. It is capable to simulate all three aforementioned benchmarks, and can be extended to support other scheduling scenarios if necessary.

F ADDITIONAL EXPERIMENTAL RESULTS

F.1 ADDITIONAL RESULTS ON PEGASUS BENCHMARK

Table 5: Experimental results on LIGO, Pegasus benchmark, with 4 different problem instance sizes (averagely 100, 200, 300, and 400 task nodes).

Method	LIGO-100			L	LIGO-200			Ll	GO-30	0	LIGO-400		
	makespan	gap/%	speed up	makespan	gap/%	speed up	mak	espan	gap/%	speed up	makespan	gap/%	speed up
GAA-PtrNet-SA GAA-PtrNet-GAT	211.0 211.8	-2.98 -2.62	2.49 2.48	460.7 460.6	-0.48 -0.50	2.49 2.50		66.9 6 6.1	-0.51 -0.63	2.50 2.51	956.9 956.6	-0.25 -0.28	2.50 2.50
PtrNet-LSTM PtrNet-CE	216.0 216.4	-0.70 -0.51	2.44 2.23	466.6 467.1	0.81 0.89	2.47 2.47	1	70.7 71.0	0.06 0.10	2.49 2.49	960.9 961.3	0.17 0.21	2.49 2.49
HEFT (heuristic baseline) Jeon et al. (2023) POMO-DAG EGS	217.5 217.0 216.7 214.2	-0.23 -0.37 -1.52	2.42 2.42 2.43 2.46	462.9 465.1 473.8 462.5	0.47 2.35 -0.09	2.49 2.48 2.43 2.49	67	70.3 70.8 75.3 69.9	0.07 0.75 -0.06	2.49 2.49 2.47 2.49	959.3 962.8 969.1 956.9	0.36 1.02 -0.25	2.49 2.49 2.47 2.50

Table 6: Experimental results on GENOME, Pegasus benchmark, with 4 different problem instance sizes (averagely 100, 200, 300, and 400 task nodes).

Method	GENOME-100			GENOME-200			GEN	NOME-	300	GENOME-400			
	makespan	gap/%	speed up	makespan	gap/%	speed up	makespan	gap/%	speed up	makespan	gap/%	speed up	
GAA-PtrNet-SA GAA-PtrNet-GAT	2435.5 2438.4	-4.61 -4.49	2.44 2.44	2351.2 2348.1	-0.94 -1.07	2.46 2.47	4723.3 4721.4	-0.60 -0.64	2.48 2.48	3451.1 3473.5	-0.06 0.559	2.48 2.47	
PtrNet-LSTM PtrNet-CE	2497.4 2493.9	-2.18 -2.32	2.38 2.38	2386.9 2386.2	0.57 -0.54	2.43 2.44	4788.4 4796.1	0.77 0.93	2.45 2.44	3505.5 3500.6	1.51 1.39	2.44 2.45	
HEFT (heuristic baseline) Jeon et al. (2023) POMO-DAG EGS	2553.1 2511.4 2472.0 2475.6	-1.63 -3.18 -3.04	2.33 2.37 2.41 2.40	2373.4 2369.9 2367.2 2356.2	-0.15 -0.26 -0.72	2.44 2.48 2.45 2.46	4751.9 4755.3 4783.2 4730.0	0.07 0.66 -0.46	2.46 2.46 2.45 2.48	3453.2 3483.2 3527.6 3453.0	0.87 2.16 -0.01	2.48 2.46 2.43 2.48	

⁴https://github.com/yd-kwon/POMO

⁵https://github.com/binqi-sun/egs

⁶https://simpy.readthedocs.io/en/

⁷https://github.com/openai/gym

F.2 ADDITIONAL RESULTS ON RANDOM GENERATED DAG WORKFLOWS

On the randomly generated DAG workflows, the performance differences between methods become more pronounced. Therefore, we present the results using bar charts. In each chart in Table 5, 6, 7 and 8, we show the outcomes when varying a single DAG generation parameter. The x-axis denoting the parameter values and each bar representing a different method. Since the makespan varies greatly across different parameter settings, the y-axis reports the speedup instead of the makespan.

It can be noticed that, although the difference is minor in most cases, in some special cases of randomly generated DAGs, the GAA-PtrNet-SA-based scheduling method obviously outperforms GAA-PtrNet-GAT. We attribute this to its attention computation mechanism (originally proposed by (Ying et al., 2021)), which does not limit attention to adjacent edges but instead leverages self-attention combined with positional encoding to globally compute attention among all nodes in the scheduling problem. This enables the model to capture dependencies even between nodes belonging to disjoint sub-DAGs, thereby enhancing its performance in such cases.

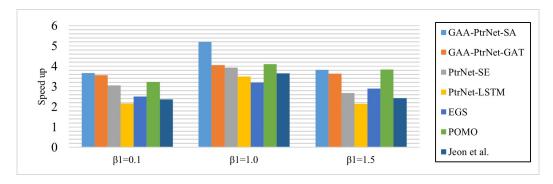


Figure 5: The results when adjusting parameter β_1 .

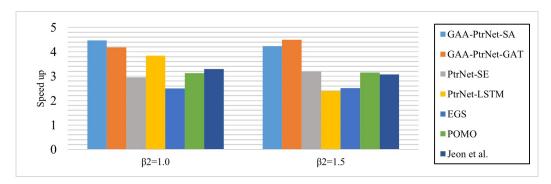


Figure 6: The results when adjusting parameter β_2 .

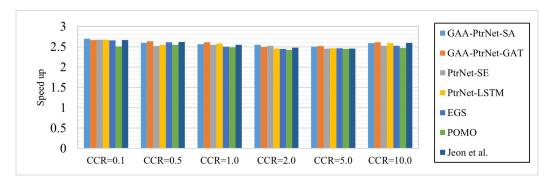


Figure 7: The results when adjusting parameter CCR.

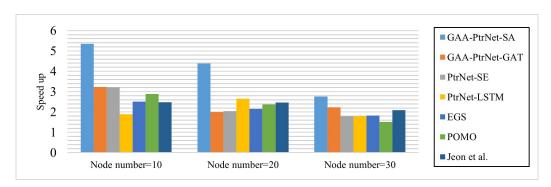


Figure 8: The results when adjusting the number of nodes in each sub-DAG.

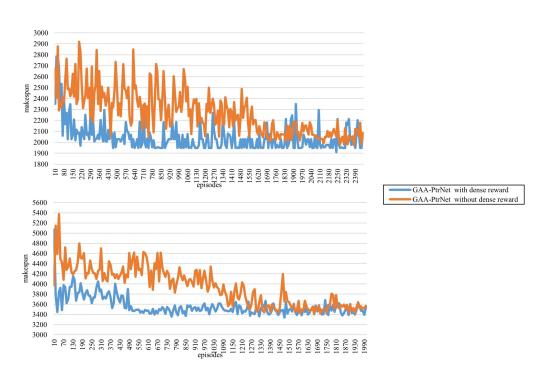


Figure 9: The curves of the makespan values evaluated on SIPHT-100 (the upper figure) and SIPHT-200 (the lower figure) during training of GAA-PtrNet, with and without dense reward signals.

G DESCRIPTION AND ANALYSIS OF THE DENSE REWARD SIGNAL

The node-level reward signal for each scheduled node is computed based on its distance to the final cost. Specifically:

- 1. Given a DAG scheduling problem G, we conduct the heuristic advantage baseline algorithm on G, obtaining each task node's individual baseline cost $C^{heur}(v_{\pi(t)})$.
- 2. For a sampled solution Solution(G), we simulate it using a SimPy-based simulator. and obtain the overall makespan C(Solution(G)) in practice. For each task node $v_{\pi(t)}$, we obtain its individual cost $C(v_{\pi(t)})$ (e.g., its finish time) through simulation.
- 3. We obtain the return-like dense reward signal R_t of each task node $v_{\pi(t)}$ by comparing the global objective C(Solution(G)) with individual cost $C(v_{\pi(t)})$, according to Equation equation 9.
- 4. To derive the advantage-like feedback A_t , we further substract each baseline from R_t , according to Equation equation 10.

 Figure 9 shows the curves of the makespan values evaluated on SIPHT-100 and SIPHT-200 during training of GAA-PtrNet, with and without dense reward signals, as a function of training epochs. It can be observed that introducing dense reward signals does not affect the quality of the final convergence, but instead leads to faster and more stable convergence.

We found that the selection of advantage baseline used in the dense reward is minor. This is because the heuristic algorithm provide constant estimates for each DAG scheduling problem instance, ensuring the advantage estimation is unbiased. Additionally, these heuristics are near-optimal in many cases, leading to similar schedules. As a result, the variance reduction benefit is preserved, while introducing little bias.

By introducing the dense reward signal, our method improves the interpretability of one-shot scheduling approaches to some extent. Specifically, the nodewise decisions, sampling probabilities and dense reward signals can be regarded as an entire MDP sampling trajectory in RL (like a trajectory by Monte Carlo sampling). This makes the interpretability of one-shot learning closer to MDP-based incremental approaches than those one-shot methods with sparse rewards.

As for the demonstration learning, it serves as the approach for the initialization of the training. Without demonstrations, we observe that the model may fall into blind exploration and even fail to converge at all. We used the genetic algorithm for DAG scheduling proposed by Zhu et al. (2016) to generate demonstrations.

H LLM USAGE IN THIS PAPER

In this paper, large language models are used only for writing embellishment and polishing, mainly focusing on certain sentences in the introduction and abstract. All the innovative points are original and no LLM was used for this.