# NN-FORMER: RETHINKING GRAPH STRUCTURE IN NEURAL ARCHITECTURE REPRESENTATION

Anonymous authors

Paper under double-blind review

#### ABSTRACT

The growing use of deep learning necessitates efficient network design and deployment, making neural predictors vital for estimating attributes such as accuracy and latency. Recently, Graph Neural Networks (GNNs) and transformers have shown promising performance in representing neural architectures. However, each method has its disadvantages. GNNs lack the capabilities to represent complicated features, while transformers face poor generalization when the depth of architecture grows. To mitigate the above problems, we rethink neural architecture topology and show that sibling nodes are pivotal while overlooked in previous research. Thus we propose a novel predictor leveraging the strengths of GNNs and transformers to learn the enhanced topology. We introduce a novel token mixer that considers siblings, and a new channel mixer named bidirectional graph isomorphism feed-forward network. Our approach consistently achieves promising performance in both accuracy and latency prediction, providing valuable insights for learning Directed Acyclic Graph (DAG) topology. The code will be released.

024 025

026 027

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

#### 1 INTRODUCTION

028 Deep neural networks have demonstrated remarkable success across various applications, high-029 lighting the significance of neural architecture design. Designing neural architectures can be quite resource-intensive. Evaluating the performance of a model necessitates training on large datasets. Measuring its inference latency and throughput involves multiple steps such as compilation, de-031 ployment, inference, and latency evaluation on various hardware platforms, incurring substantial human effort and resources. One strategy to mitigate these challenges is to predict network attributes 033 with machine learning predictors. By feeding the network structure and hyperparameters into these 034 predictors, valuable attributes of the network can be estimated with just a single feedforward pass, e.g., accuracy on a validation set or inference times on specific hardware configurations. This predictive approach has been successfully applied in various tasks including neural architecture search (Xu 037 et al., 2021; Luo et al., 2018; Wen et al., 2020; Lu et al., 2021; Yi et al., 2023; 2024) and hardware 038 deployment (Zhang et al., 2021; Kaufman et al., 2021; Dudziak et al., 2020; Liu et al., 2022; Yi et al., 2023; 2024), yielding promising outcomes in improving the efficiency and effectiveness of network design. 040

With the recent development of transformers, various transformer-based frameworks have been introduced (Lu et al., 2021; Yi et al., 2023; 2024). Transformers have strengths in global modeling and

<sup>041</sup> Previous neural predictors model the neural architecture as a Directed Acyclic Graph (DAG) (Wen 042 et al., 2020; Li et al., 2020; Shi et al., 2020; Dudziak et al., 2020; Liu et al., 2022; Dong et al., 043 2022; Luo et al., 2023) and utilize Graph Neural Networks (GNNs) or Transformers to extract 044 neural architecture representation. GNNs have emerged as an intuitive solution for learning graph representations (Wen et al., 2020; Li et al., 2020; Shi et al., 2020; Dudziak et al., 2020; Liu et al., 2022), which leverage the graph Laplacian and integrate adjacency information to learn the graph 046 topology. GNN-based predictors show strong generalization ability, yet their performance may 047 not be optimal. This is attributed to the structural bias in the message-passing mechanism, which 048 typically relies solely on adjacency information. As illustrated in Figure. 1(a) and (b), GCNs (Kipf & Welling, 2016) aggregate the forward and backward adjacent nodes without discrimination, and GATs (Veličković et al., 2018) aggregate them with dynamic weights. Both of them are limited to 051 adjacent information. 052



067 Figure 1: Comparison of different methods on DAG representation of neural architectures. (a) GCNs (Kipf & Welling, 2016) aggregate adjacent information without discrimination. (b) 068 GATs (Veličković et al., 2018) distinguish adjacent operations, while are still constrained to adjacent 069 nodes. (c) Vanilla transformers (Vaswani et al., 2017) aggregate weighted global information, which can result in poor generalization as the network depth increases. (d) Transformers on directed transi-071 tive closure (Dong et al., 2022; Luo et al., 2023) aggregate the successor information but still suffer 072 from poor generalization. (e) Our method aggregates sibling information with weighted coefficients. 073 Sibling nodes could extract complementary features in accuracy prediction and allow for concurrent 074 execution in latency prediction. 075

dynamic weight adjustments, hence could extract strong features. Despite the promising performance, they still exhibit several shortcomings. One particular challenge of transformer is related to the 079 long-range receptive field, as depicted in Figure. 1(c) and (d), which can lead to poor generalization performance on deep architectures (Yi et al., 2023; 2024). The vanilla transformer (Vaswani et al., 081 2017; Lu et al., 2021; Yi et al., 2023) have a global receptive field, and recent studies proposed 082 transformers on directed transitive closure (Dong et al., 2022; Luo et al., 2023). Both methods conduct 083 long-range attention that could mix up the information from operations far away, especially when the 084 depth of the input architecture increases to hundreds of layers. For example, NAR-Former (Yi et al., 085 2023; 2024) has illustrated that transformer predictors with global attention struggle in deep network latency prediction, leading to worse performance than GNNs (Liu et al., 2022; Yi et al., 2024).

087 To study a more effective neural predictor, we rethink the DAG topology and show that the commonly 088 used topological information is not suitable for the neural architecture representation. Most of recent works focus on modeling the relationship of preceding and succeeding operations (Dong et al., 2022; 090 Luo et al., 2023). However, it is essential to recognize the importance of "sibling nodes", which share 091 a common parent or child node with the current node as shown in Figure. 1(e). They often exhibit 092 strong connections to the current nodes in neural architecture representation. For example, in the accuracy prediction task, parallel branches may extract complementary features, hence enhancing overall model performance. Furthermore, operations that share the same parent or child node can 094 be executed simultaneously, potentially reducing inference latency. On the contrary, long-range 095 dependency might not be crucial, given that features typically propagate node-by-node within the 096 architecture. Previous methods have not explicitly leveraged sibling cues.

Based on the analysis above, we introduce a new model for neural architecture representation, named
Neural Network transFormer (NN-Former). It leverages the strengths of GNNs and transformers,
exhibiting good generalization and high performance. For the token mixing module, we utilize a
self-attention mechanism of transformers to extract dynamic weights for capturing complex features.
We explicitly learn the adjacency and sibling nodes' features to enhance the topological information.
For the channel mixing module, we use a bidirectional graph isomorphism feedforward network. It
learns strong graph topology information such that the position encoding is no longer necessary.

Extensive experiments reveal that 1) our approach surpasses existing methods in both accuracy
 prediction and latency prediction, demonstrating expressivity and generalization ability, and 2) our
 method has good scalability on both cell-structured architectures and complete neural networks that
 have hundreds of operations. To the best of our knowledge, this is an original work that leverages

sibling cues in neural predictor. Integrating strengths of both GNNs and transformers effectively
 guarantees its promising performance in. The importance of sibling nodes also provides a valuable
 insight into rethinking DAG topology representation in future research.

111 112

113

# 2 RELATED WORKS

114 Neural Architecture Representation Learning. Neural architecture representation learning es-115 timates network attributes without actual training or deployment, resulting in significant resource 116 savings. Accuracy predictors forecast the evaluation accuracy, avoiding the resource-intensive process 117 of network training in neural architecture search (Liu et al., 2018; White et al., 2021; Deng et al., 118 2017; Luo et al., 2018; 2020; Cai et al., 2019; Zhang et al., 2018; Li et al., 2020; Shi et al., 2020; 119 Chen et al., 2021; Yan et al., 2020; Lu et al., 2021; Yi et al., 2023; 2024). Additionally, latency 120 prediction can estimate the inference latency without actual deployment, saving time and materials 121 for engineering application (Dudziak et al., 2020; Zhang et al., 2021; Kaufman et al., 2021; Liu et al., 2022). Given the complex connections between operations and the one-way message-passing 122 mechanism, the neural network is better described as a DAG, with the connection between nodes 123 represented by the adjacency matrix. Consequently, graph-based (Zhang et al., 2018; Li et al., 2020; 124 Shi et al., 2020; Chen et al., 2021; Yan et al., 2020) and transformer-based (Lu et al., 2021; Yi 125 et al., 2023; 2024) predictors have been employed to learn the representation of neural architectures. 126 Both methods achieve promising results in neural architecture representation, while they still face 127 challenges. In this paper, we absorb the strengths of both methods and delve into the topological 128 relationship. 129

130 Message-Passing Graph Neural Networks. Most contemporary graph neural networks can be 131 expressed within the framework of the message-passing architecture (Gilmer et al., 2017; Kipf & Welling, 2016; Hamilton et al., 2017; Veličković et al., 2018; Xu et al., 2018; You et al., 2020). In 132 this framework, node representations are computed iteratively by aggregating the embeddings of 133 their neighboring nodes, and a final graph representation can be obtained by aggregating the node 134 embeddings, such as GCN (Kipf & Welling, 2016), GAT (Veličković et al., 2018), GIN (Xu et al., 135 2018), etc. GNN-based models have emerged as a prominent and widely adopted approach for 136 neural network representation learning (Zhang et al., 2018; Li et al., 2020; Shi et al., 2020; Chen 137 et al., 2021; Yan et al., 2020). The straightforward structural characteristics of GNNs contribute to 138 strong generalization ability, yet also necessitate further improvement in the performance. Enhancing 139 topological information and dynamic bi-directional aggregation is a promising approach. 140

141 Transformers on Graphs. Recently, transformer has been introduced into graph representation 142 learning (Dwivedi & Bresson, 2020; Wu et al., 2021; Dong et al., 2022; Luo et al., 2023), together with 143 network architecture representation learning (Lu et al., 2021; Yi et al., 2023; 2024). TNASP (Lu et al., 144 2021) inputs the sum of the operation type embedding matrix and Laplacian matrix into the standard 145 transformer. NAR-Former (Yi et al., 2023) encodes each operation and connection information into a token and inputs all tokens into a proposed multi-stage fusion transformer. NAR-Former V2 (Yi 146 et al., 2024) introduced a graph-aided transformer block, which can handle both cell-structured 147 networks and entire networks. However, transformers face challenges of poor generalization when 148 the network goes deeper, with global attention mixing up the far away information (Yi et al., 2024). 149 To address this limitation, we propose a novel predictor that harnesses the strengths of both GNNs 150 and transformers, allowing it to extract both topology features and dynamic weights. This approach 151 enhances the model's ability to extract valuable insights and maintains good generalization. 152

152

153 **Neural Networks over DAGs.** The inductive bias inherent in DAGs has led to specialized neural 154 predictors. GNNs designed for DAGs typically compute graph embeddings using a message-passing 155 framework (Thost & Chen, 2021). On the other hand, transformers applied to DAGs often incorporate 156 the depth of nodes (Kotnis et al., 2021; Luo et al., 2023) or Laplacian (Gagrani et al., 2022) as the 157 position encoding, which may seem non-intuitive for integrating structural information into the model. 158 Additionally, transformer-based models frequently use transitive closure (Dong et al., 2022; Luo 159 et al., 2023) as attention masks, leading to poor generalization as mentioned above. Some hybrid methods with GNNs and Transformers are not tailored to neural architecture representation and also 160 face similar challenges (Ying et al., 2021; Wu et al., 2021). This paper proposes a novel hybrid model 161 with enhanced topological information from sibling nodes.

177

178

179

180 181

182 183

184

213

214



Figure 2: The proposed NN-Former framework. We introduce adjacency and sibling attention masks in the Adjacency-Sibling Multihead Attention (ASMA) to learn graph topology information. We also introduce adjacency aggregation in the Bidirectional Graph Isomorphism Feed-Forward Network (BGIFFN) to enhance the topology structure.

#### 3 METHODS

#### 3.1 OVERVIEW

185 We adopt a commonly used graph representation of neural architectures (Lu et al., 2021; Yi et al., 2023; 2024; Dong et al., 2022; Luo et al., 2023; Liu et al., 2022). An architecture with n operations 187 is referred to as a Graph G = (V, E, Z), with node set V, edge set  $E \subseteq V \times V$ , and node features  $Z \in \mathbb{R}^{n \times d}$ . Each operation is denoted as a node in V such that |V| = n. The edge set E is often given in form of an adjacency matrix  $A \in \{0,1\}^{n \times n}$ , where  $A_{ij} = 1$  denotes a directed edge 188 189 190 from node i to node j. Each row of Z represents the feature vector of one node, *i.e.*, operation type and hyperparameters, with the number of nodes n and feature dimension d. Unlike previous 191 methods (Lu et al., 2021; Yi et al., 2023), our predictor is strong and position encoding is unnecessary. 192 For simplicity, Z is encoded with one-hot encoding for operation type and sinusoidal encoding for 193 operation attributes as (Yi et al., 2024). The neural architecture representation (Luo et al., 2018; Wen 194 et al., 2020; Xu et al., 2021; Lu et al., 2021; Yi et al., 2023; 2024; Liu et al., 2022) utilizes a predictor 195  $f_{\theta}(\cdot)$  with parameters  $\theta$  to estimate specific attributes of candidate architectures, e.g., validation 196 accuracy or inference latency: 197

$$\hat{y} = f_{\theta} \left( \boldsymbol{Z}, \boldsymbol{A} \right), \tag{1}$$

where  $\hat{y}$  denotes the predicted attribute of the architecture. 199

As illustrated in Figure. 2, our approach uses a transformer as the baseline and incorporates dis-200 criminative topological features to pursue a strong predictor. Previous transformer-based predictors 201 considered adjacent propagation (Lu et al., 2021; Yi et al., 2023; 2024) or transitive closure (Dong 202 et al., 2022; Luo et al., 2023) as the graph structure information. Global attention is effective in 203 shallow network prediction as shown in previous works such as TNASP (Lu et al., 2021) and NAR-204 Former (Yi et al., 2023). However, with the network depth increasing, there is a decrease in the 205 generalization of global attention as shown in Yi et al. (2024). Global attention may be biased 206 towards training data and demonstrate poor generalization performance. To build a general neural 207 predictor for the range of all depths, we propose a non-global neural predictor that outperforms 208 the previous methods on both accuracy and latency predictions. As we discussed in Section 1, 209 sibling nodes have a strong relationship with the current nodes and also provide useful information 210 in accuracy and latency prediction. Thus we introduce an Adjacency-Sibling Multi-head Attention (ASMA) in the self-attention layer to learn the local features. As the sibling relationship can be 211 calculated from adjacency matrix A, our ASMA is formulated as: 212

> $\hat{\boldsymbol{H}}^{l-1} = \text{ASMA}\left(\text{LN}\left(\boldsymbol{H}^{l-1}\right), \boldsymbol{A}\right) + \boldsymbol{H}^{l-1},$ (2)

where  $H^{l}$  denotes the feature for the layer l and LN denotes layer normalization. ASMA injects 215 topological information into the transformer, thereby augmenting the capability of Directed Acyclic

222

227

228 229

230

242 243

259

260 261

262

263

264

265

266

Graph (DAG) representation learning. In the channel-mixing part, we introduce a Bidirectional Graph
 Isomorphism Feed-Forward Network (BGIFFN). This module extracts strong topology features and alleviates the necessity of complex position encoding:

$$\boldsymbol{H}^{l} = \text{BGIFFN}\left(\text{LN}\left(\hat{\boldsymbol{H}}^{l-1}\right), \boldsymbol{A}\right) + \hat{\boldsymbol{H}}^{l-1}.$$
(3)

As for input and output, the first layer feature  $H^0$  and the last layer feature  $H^L$  are related to the input and output in the following way:

$$\boldsymbol{H}^{0} = \mathrm{LN}\left(\mathrm{FC}\left(\boldsymbol{Z}\right)\right),\tag{4}$$

$$\hat{y} = \text{FC}\left(\text{ReLU}\left(\text{FC}\left(\boldsymbol{H}^{L}\right)\right)\right),\tag{5}$$

where FC denotes fully-connected layer. The following parts proceed to introduce ASMA and BGIFFN in detail.

#### 3.2 ADJACENCY-SIBLING MULTIHEAD ATTENTION

Given a node, we define its sibling nodes as those that share the same parents or children. To identify these sibling nodes, we use the adjacency matrix A and its transpose  $A^T$ . Specifically, nodes sharing the same parent nodes are indicated by the non-zero positions in the matrix product  $AA^T$ , reflecting the backward mapping followed by the forward mapping. Similarly, nodes sharing the same children nodes are identified through the matrix product  $A^TA$ . In this way, we can identify whether there is a sibling relationship between each pair of nodes.

To inject topological information, we introduce a novel multi-head attention module. As shown in Figure. 2, we use four-head attention, where each head uses an attention mask indicating a specific topology. These masks include forward adjacency A, backward adjacency  $A^T$ , siblings with the same parents  $AA^T$ , and siblings with the same children  $A^TA$ , respectively. The proposed ASMA is denoted as:

$$ASMA(\boldsymbol{H}) = Concat(\boldsymbol{X}_1, \boldsymbol{X}_2, \boldsymbol{X}_3, \boldsymbol{X}_4) \boldsymbol{W}^O,$$
(6)

$$\boldsymbol{X}_{1} = \sigma\left(\left(\boldsymbol{Q}_{1}\boldsymbol{K}_{1}^{T}\circ(\boldsymbol{I}+\boldsymbol{A})\right)/\sqrt{h}\right)\boldsymbol{V}_{1},\tag{7}$$

$$\boldsymbol{X}_{2} = \sigma \left( \left( \boldsymbol{Q}_{2} \boldsymbol{K}_{2}^{T} \circ \left( \boldsymbol{I} + \boldsymbol{A}^{T} \right) \right) / \sqrt{h} \right) \boldsymbol{V}_{2}, \tag{8}$$

$$\boldsymbol{X}_{3} = \sigma \left( \left( \boldsymbol{Q}_{3} \boldsymbol{K}_{3}^{T} \circ \left( \boldsymbol{I} + \boldsymbol{A} \boldsymbol{A}^{T} \right) \right) / \sqrt{h} \right) \boldsymbol{V}_{3}, \tag{9}$$

$$\boldsymbol{X}_{4} = \sigma \left( \left( \boldsymbol{Q}_{4} \boldsymbol{K}_{4}^{T} \circ \left( \boldsymbol{I} + \boldsymbol{A}^{T} \boldsymbol{A} \right) \right) / \sqrt{h} \right) \boldsymbol{V}_{4}, \tag{10}$$

where  $X_i$  denote the *i*-th head feature and  $Q_i = HW_i^Q$ ,  $K_i = HW_i^K$ ,  $V_i = HW_i^V$  denotes 251 the query, key, and value for each head, respectively.  $\sigma$  is the softmax operation, and h denotes the 252 dimension of each head. I is introduced to contain self-position information.  $\circ$  is an elementwise 253 masking operation, which constrains the attention to the non-zero positions of the mask matrix. 254 For example, the first head  $H_1$  utilizes an attention mask of I + A, which means it only conducts 255 attention on the self-position and forward adjacency position. ASMA decouples the local structure 256 information into 4 different perspectives. This module extracts diverse topological information and 257 thus exhibits enhanced representation power in modeling neural architecture. 258

#### 3.3 BIDIRECTIONAL GRAPH ISOMORPHISM FEED-FORWARD NETWORK

To further enhance the topology information, we propose a bidirectional graph isomorphism feedforward network. We utilize the adjacency matrix A and its transpose  $A^T$  to aggregate the forward and backward adjacency positions in the feedforward module. The BGIFFN is formulated as:

BGIFFN 
$$(\boldsymbol{H}, \boldsymbol{A}) = \text{ReLU} (\boldsymbol{H}\boldsymbol{W}_1 + \boldsymbol{H}_q) \boldsymbol{W}_2,$$
 (11)

$$\boldsymbol{H}_{g} = \operatorname{Concat}\left(\operatorname{GC}\left(\boldsymbol{H},\boldsymbol{A}\right),\operatorname{GC}\left(\boldsymbol{H},\boldsymbol{A}^{T}\right)\right),\tag{12}$$

where  $H_g$  denotes the output features of the graph convolution and  $W_1$ ,  $W_2$  denote the parameters of linear transformation. GC denotes the graph convolution, which is a simplified form of GCN(Kipf & Welling, 2016):

$$GC(\boldsymbol{H}, \boldsymbol{A}) = \boldsymbol{A}\boldsymbol{H}\boldsymbol{W},\tag{13}$$

Backbone	Method	Publication	0.02% (100)	Training S 0.04% (172)	amples 0.1% (424)	1% (4236)
CNN	ReNAS (Xu et al., 2021)	CVPR 2021	-	-	0.657	0.816
LSTM	NAO (Luo et al., 2018)	NeurIPS 2018	0.501	0.566	0.666	0.775
GNN	NP (Wen et al., 2020) GATES (Ning et al., 2020) GMAE-NAS (Jing et al., 2022)	ECCV 2020 ECCV 2020 IJCAI 2022	0.391 0.605 0.666	0.545 0.659 0.697	0.679 0.691 0.732	0.769 0.822 0.775
Transformer	Graphormer (Ying et al., 2021) TNASP (Lu et al., 2021) NAR-Former (Yi et al., 2023) PINAT (Lu et al., 2023)	NeurIPS 2021 NeurIPS 2021 CVPR 2023 AAAI 2024	0.564 0.600 0.632 0.679	0.580 0.669 0.653 0.715	0.611 0.705 0.765 0.772	0.797 0.820 0.871 0.846
Hybrid	GraphTrans (Wu et al., 2021) NAR-Former V2 (Yi et al., 2024) NN-Former (Ours)	NeurIPS 2021 NeurIPS 2023 -	0.330 0.663 <b>0.709</b>	0.472 0.704 <b>0.765</b>	0.602 0.773 <b>0.809</b>	0.700 0.861 <b>0.877</b>

Table 1: Accuracy prediction results on NAS-Bench-101 (Ying et al., 2019). We use different proportions of data as the training set and report Kendall's Tau on the whole dataset.

where W denotes the parameters of fully-connected layer. Note that we use the directed adjacency matrix rather than graph Laplacian, which makes it simpler and stronger. With GC (H, A) and GC  $(H, A^T)$ , we obtain the forward features and backward features. Since we concatenate the two directional features, the BGIFFN will learn forward propagation in one half of the channels, and backward propagation in the other. We will show that BGIFFN demonstrates bidirectional graph isomorphism in Appendix A.3, which enhances topology information.

291 292 293

294 295

296

297

298

287

288

289

290

### 4 EXPERIMENTS

We conduct experiments on two tasks, namely accuracy prediction and latency prediction. For accuracy prediction, we evaluate the ranking performance of NN-Former on two benchmarks NAS-Bench-101 (Ying et al., 2019) and NAS-Bench-201 (Dong & Yang, 2020). For latency performance, we conduct experiments on NNLQ (Liu et al., 2022). A series of ablation experiments are conducted to demonstrate the effectiveness of our design. More details will be provided in the appendix.

#### 4.1 ACCURACY PREDICTION

We conduct accuracy prediction on NAS-Bench-101 (Ying et al., 2019) and NAS-Bench-201 (Dong & Yang, 2020). Both datasets adopt cell-structured architectures. The NAS-Bench-101 (Ying et al., 2019) dataset contains 423,624 unique architectures, each comprising 9 repeated cells with a maximum of 7 nodes and 9 edges per cell. Similar to the NAS-Bench-101, the architectures in NAS-Bench-201 (Dong & Yang, 2020) are also built using repeated cells. It presents 15,625 distinct cell candidates, each composed of 4 nodes and 6 edges. We report Kendall's Tau as the previous methods (Lu et al., 2021; Ning et al., 2020; Yi et al., 2023).

Experiments on NAS-Bench-101. We implemented the configuration outlined in TNASP (Lu et al., 2021) to train our predictor on subsets of 0.02%, 0.04%, 0.1%, and 1% of the entire dataset. Subsequently, we utilized the complete dataset as the test set and computed Kendall's Tau to evaluate the performance. The results are detailed in Table 1. Our predictor consistently outperforms baseline methods, such as CNNs, LSTMs, GNNs, Transformers, and hybrid GNNs and Transformers. This result underscores the superior predictive capability of NN-Former in determining neural architecture performance.

Experiments on NAS-Bench-201. We employ a comparable experimental setup to NAS-Bench-101,
 *i.e.*, training predictors on different subsets of 1%, 3%, 5%, and 10%, and then evaluating them on
 the complete dataset. The results are depicted in Table 2. NN-Former surpasses other methods in
 all scenarios except for the 10% subsets. Note that our method aims at unified prediction for both
 accuracy and latency, it is acceptable that our method achieves comparable results. We outperform
 NAR-Former V2 (Yi et al., 2024) for all setups, which has a similar unifying motivation. Additionally,
 neural architecture search prefers high generalization performance with fewer training samples,
 resulting in significant resource savings. More details are discussed in Section B.1.2.

325 326 327

339

355

362

324

Table 2: Accuracy prediction results on NAS-Bench-201 (Dong & Yang, 2020). We use different proportions of data as the training set and report Kendall's Tau on the whole dataset.

Backbone	Method		Publication		1% (156)	Trainin 3% (469)	g Samples 5% (781)	10% (1563)
LSTM	NAO (Luo et al., 2018)		NeurIPS 2018		0.493	0.470	0.522	0.526
GNN	NP (Wen et al., 2020)	1	ECCV 2020		0.413	0.584	0.634	0.646
Transformer	Graphormer (Ying et al., 2021) TNASP (Lu et al., 2021) NAR-Former (Yi et al., 2023) PINAT (Lu et al., 2023)		NeurIPS 2021 NeurIPS 2021 CVPR 2023 AAAI 2024		0.630 0.589 0.660 0.631	0.680 0.640 0.790 0.706	0.719 0.689 0.849 0.761	0.776 0.724 <b>0.901</b> 0.784
Hybrid	GraphTrans (Wu et al., 2021) NAR-Former V2 (Yi et al., 2024) NN-Former (Ours)		NeurIPS 2021 NeurIPS 2023		0.409 0.752 <b>0.804</b>	0.550 0.846 <b>0.860</b>	0.588 0.874 <b>0.879</b>	0.673 0.888 0.890

Table 3: In domain latency prediction on NNLQ (Liu et al., 2022). Training and test on the same distribution.

Test Model	NNLP avg / best	MAPE↓ NAR-Former V2 avg / best	Ours avg / best	NNLP avg / best	Acc(10%)↑ NAR-Former V2 avg / best	Ours avg / bes
All	3.47 / 3.44	3.07 / 3.00	2.85 / 2.65	95.25 / 95.50	96.41 / 96.30	97.45 / 97
AlexNet	6.37/6.21	6.18 / 5.97	4.69 / 4.61	81.75 / 84.50	81.90 / 84.00	90.50/91
EfficientNet	3.04 / 2.82	2.34 / 2.22	2.31 / 2.21	98.00 / 97.00	98.50 / 100.0	99.00 / 10
GoogleNet	4.18/4.12	3.63 / 3.46	3.48 / 3.39	93.70/93.50	95.95 / 95.50	97.15/97
MnasNet	2.60 / 2.46	1.80 / 1.70	1.52 / 1.48	97.70/98.50	99.70 / 100.0	99.50 / 10
MobileNetV2	2.47/2.37	1.83 / 1.72	1.54 / 1.50	99.30 / 99.50	99.90 / 100.0	99.60 / 10
MobileNetV3	3.50/3.43	3.12/2.98	3.17 / 2.99	95.35 / 96.00	96.75 / 98.00	96.50/97
NasBench201	1.46/1.31	1.82/1.18	1.11/0.96	100.0 / 100.0	100.0 / 100.0	100.0 / 10
SqueezeNet	4.03 / 3.97	3.54 / 3.34	3.09 / 3.08	93.25 / 93.00	95.95 / 96.50	97.70/98
VĠG	3.73/3.63	3.51/3.29	2.94 / 2.89	95.25 / 96.50	95.85 / 96.00	95.80/96
ResNet	3.34/3.25	3.11/2.89	2.66 / 2.47	98.40/98.50	98.55 / 99.00	99.45 / 99

4.2 LATENCY PREDICTION

We employ NNLQ as the latency prediction task. NNLQ (Liu et al., 2022) includes 20,000 deeplearning networks and their respective latencies on the specified hardware. This dataset encompasses 10 distinct network types, with 2,000 networks for each type. The depth of each architecture varies from tens to hundreds of operations, requiring the scalability of the neural predictor. In line with NNLP, the Mean Absolute Percentage Error (MAPE) and Error Bound Accuracy (Acc( $\delta$ )) are employed to assess the disparities between latency predictions and actual values.

We conduct the experiments on two scenarios following (Yi et al., 2024). In the first in-domain scenario, the training and testing sets are from the same distribution. The results are shown in Table 3. When testing with all test samples, the average MAPE of our methods is 0.62% lower than the NNLP (Liu et al., 2022) and 0.22% lower than the NAR-Former V2 (Yi et al., 2024). The average Acc(10%) is 2.20% higher than the NNLP and 1.04% higher than the NAR-Former V2. When tested on various types of network data separately, previous methods fail on specific model types, especially on AlexNet, while our method largely mitigates this challenge and obtains a balanced performance on each model type.

370 The other out-of-domain scenario is more significant, as it involves inferring an unseen network type 371 during the evaluation. The results in Table 4 indicate that relying solely on FLOPs and memory access 372 data is insufficient for predicting latency. Due to the disparity between kernel delay accumulation and 373 actual latency, kernel-based approaches such as nn-Meter (Zhang et al., 2021) and TPU (Kaufman 374 et al., 2021) exhibit inferior performance compared to GNNs (NNLP (Liu et al., 2022)) and hybrid 375 models (NAR-Former V2 (Yi et al., 2024)). Leveraging enriched topological information, our method achieves the highest MAPE and Acc(10%) among the average metrics of the ten experimental sets. 376 In comparison to the runner-up NAR-Former V2 (Yi et al., 2024), our approach demonstrates a 377 substantial 11.61% increase in average Acc(10%).

Table 4: Out of domain latency prediction on NNLQ (Liu et al., 2022). "Test Model = AlexNet" means that only AlexNet models are used for testing, and the other 9 model families are used for training. The best results refer to the lowest MAPE and corresponding ACC (10%) in 10 independent experiments.

Test Model	FLOPs	FLOPs +MAC	nn-Meter	TPU	BRP-NAS	NNLP (avg / best)	NAR-Former V2 (avg / best)	Ours (avg / best)
					MA	PE↓		
AlexNet	44.65	15.45	7.20	10.55	31.68	10.64 / 9.71	24.28 / 18.29	11.47 / 11.12
EfficientNet	58.36	53.96	18.93	16.74	51.97	21.46 / 18.72	13.20 / 11.37	5.13 / 4.81
GoogleNet	30.76	32.54	11.71	8.10	25.48	13.28 / 10.90	6.61 / 6.15	6.74 / 6.65
MnasNet	40.31	35.96	10.69	11.61	17.26	12.07 / 10.86	7.16 / 5.93	2.71 / 2.54
MobileNetV2	37.42	35.27	6.43	12.68	20.42	8.87 / 7.34	6.73 / 5.65	4.17 / 3.66
MobileNetV3	64.64	57.13	35.27	9.97	58.13	14.57 / 13.17	9.06 / 8.72	9.07 / 9.03
NasBench201	80.41	33.52	9.57	58.94	13.28	9.60 / 8.19	9.21 / 7.89	7.93 / 7.71
ResNet	21.18	18.91	15.58	20.05	15.84	7.54 / 7.12	6.80 / 6.44	7.49/7.38
SqueezeNet	29.89	23.19	18.69	24.60	42.55	9.84/9.52	7.08 / 6.56	9.08 / 7.05
VGG	69.34	66.63	19.47	38.73	30.95	7.60 / 7.17	15.40 / 14.26	20.12 / 19.64
Average	47.70	37.26	15.35	21.20	30.76	11.55 / 10.27	10.55/9.13	8.39 / 7.96
					Acc(1	0%)↑		
AlexNet	6.55	40.50	75.45	57.10	15.20	59.07 / 64.40	24.65 / 28.60	56.08 / 57.1
EfficientNet	0.05	0.05	23.40	17.00	0.10	25.37 / 28.80	44.01 / 50.20	90.85 / 90.9
GoogleNet	12.75	9.80	47.40	69.00	12.55	36.30/48.75	80.10 / 83.35	80.43 / 83.4
MnasNet	6.20	9.80	60.95	44.65	34.30	55.89/61.25	73.46 / 81.60	98.65 / 98.7
MobileNetV2	6.90	8.05	80.75	33.95	29.05	63.03 / 72.50	78.45 / 83.80	94.90 / 96.8
MobileNetV3	0.05	0.05	23.45	64.25	13.85	43.26 / 49.65	68.43 / 70.50	74.18 / 74.3
NasBench201	0.00	10.55	60.65	2.50	43.45	60.70 / 70.60	63.13 / 71.70	69.90 / 71.1
ResNet	26.50	29.80	39.45	27.30	39.80	72.88 / 76.40	77.24 / 79.70	70.83 / 71.5
SqueezeNet	16.10	21.35	36.20	25.65	11.85	58.69 / 60.40	75.01 / 79.25	77.85 / 80.9
VGG	4.80	2.10	26.50	2.60	13.20	71.04 / 73.75	45.21 / 45.30	29.40 / 29.8
Average	7.99	13.20	47.42	34.40	21.34	54.62 / 60.65	62.70 / 67.40	74.31 / 75.4
-	0							

Table 5: **Ablation studies on the proposed ASMA and BGIFFN.** (a) Ablation study of ASMA on NNLQ (Liu et al., 2022). Results on the NAS-Bench-201 family are reported. (b) Ablation study of ASMA and BGIFFN on NAS-Bench-101 (Ying et al., 2019). All experiments are conducted on the 0.04% training set.

(a)			(b)	
Attention   MAPE $\downarrow$ Acc(10%) $\uparrow$	_	Attention	FFN	Kendall's Tau↑
Global         10.83%         58.45%           ASMA         7.93%         69.90%	-	Global ASMA Global ASMA	vanilla vanilla BGIFFN BGIFFN	0.4598 0.6538 0.7656 0.7654

#### 4.3 ABLATION STUDIES

In this section, we perform a series of ablation experiments on the NAS-Bench-101 (Ying et al., 2019)
 and NNLQ (Liu et al., 2022) datasets to analyze the effects of the proposed modifications. First, we
 assess the performance of ASMA and BGIFFN over a vanilla transformer baseline. Second, we verify
 the design for both modules, especially by examining the significance of the topological information.

Ablation on ASMA on latency prediction. To evaluate the generalization ability of ASMA, we conducted experiments on NNLQ (Liu et al., 2022) under the out-of-domain setting. We use NAS-Bench-201 as our target model type due to its intricate connections between operations. The results are shown in Table 5a, where we keep the number of heads unchanged ablate on the attention mask. The model with ASMA shows a large performance enhancement of 11.45% Acc(10%) in contrast to the one utilizing global attention. It indicates that global attention presents a noticeable disparity, underscoring the advantages of local features when building a unified neural predictor.

Ablation on ASMA and BGIFFN on accuracy prediction. To evaluate the effectiveness of ASMA and BGIFFN, we conducted accuracy prediction experiments using the 0.04% setting of NAS-Bench-101 as discussed in Section 4.1. The findings are detailed in Table 5b. The baseline method utilizes a vanilla transformer and the result is respectable since it does not incorporate any topology information. Introducing ASMA, which integrates adjacency and sibling information, leads to an enhancement of 0.6538. Incorporating BGIFFN leads to further enhancement of 0.7656. It indicates that the ASMA and BGIFFN modules effectively capture structural features within the models. Furthermore, when

Table 6: Ablation studies on ASMA design. All experiments are conducted on the NAS-Bench101 (Ying et al., 2019) with the 0.04% training set. (a) Ablation study on the topological structure
of ASMA. We explore different attention masks for the 4 heads. (b) Ablation study on the position
encoding with ASMA. We explore different ways of position encoding when utilizing ASMA.

(a)		(b)	
Row   Attention Mask	Kendall's Tau↑	ASMA design	Kendall's Tau↑
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	0.7522 0.7545 0.7566 0.7573 <b>0.7654</b>	Ours +NAR PE (Yi et al., 2023) +Laplacian (Lu et al., 2021)	<b>0.7654</b> 0.7449 0.7063

Table 7: **Ablation studies on BGIFFN design.** All experiments are conducted on NAS-Bench-101 (Ying et al., 2019) with the 0.04% training set. (a) Ablation study on the topological structure of BGIFFN. We explore different ways of structure aggregation in BGIFFN. (b) Ablation study on the BGIFFN design. We explore different ways of adjacency aggregation when utilizing BGIFFN.

	(a)		(b)	
Row	BGIFFN	Kendall's Tau↑	BGIFFN design	Kendall's Tau↑
1	$A_{_{TT}}$	0.7253	Ours	0.7654
2	$A_{-}^{I}$	0.7501	add→multiply	0.7076
3	$\boldsymbol{A}, \boldsymbol{A}^T, \boldsymbol{A}^T \boldsymbol{A}, \boldsymbol{A} \boldsymbol{A}^T$	0.7470	$\rightarrow$ GCN (Laplacian)	0.7296
4	$oldsymbol{A},oldsymbol{A}^T$	0.7654	$\rightarrow$ GAT	0.6973

ASMA and BGIFFN are combined, the model achieves a performance score of 0.7654 which is comparable to the global attention mechanism. However, we have highlighted the limitations of global attention in latency prediction, and our ASMA achieved the best trade-off across both accuracy and latency prediction tasks.

Ablation on topological structure of ASMA. Next, we verify the design of the ASMA and BGIFFN. First, we examine the designs of ASMA in Table 6. The performance of different attention masks is shown in Table 6a. Maintaining a consistent number of heads at 4, we modify the attention mask for each head. Rows 1 and 2 exclusively utilize forward or backward adjacency across all 4 heads. Row 3 combines forward and backward adjacency and improves performance. Row 4 investigates the impact of predecessors and successors, indicating only marginal enhancement. It shows that empirical topological information in DAG tasks (Dong et al., 2022; Luo et al., 2023) is helpless in neural architecture representation. Row 5 combines the adjacency and sibling nodes, achieving the highest performance. These results highlight the significance of sibling nodes. It also shows that the design of ASMA is robust and logically sound, showcasing its effectiveness in capturing architectural patterns. 

Ablation on Position Encoding (PE). We also explored the impact of PE on our model. Traditionally, transformers have heavily relied on PE to capture structural information. However, our approach makes this reliance unnecessary because our method inherently incorporates abundant topological information. As shown in Table 6b, we experiment with the inclusion of position encoding in our framework, *i.e.*, NAR PE tailored for neural architecture representation in NAR-Former (Yi et al., 2023), and Laplacian position encoding in TNASP (Luo et al., 2023). The results suggest that they have no improvement in the performance of NN-Former. It indicates that our method presents exceptional structural learning capabilities. 

Ablation on topological structure of BGIFFN. As illustrated in Table 7a, this experiment maintains
 the total parameters constant while adjusting the number of splits in the graph convolution branch.
 In Rows 1 and 2, a single split is retained, and the forward or backward adjacent convolution is
 conducted. Moving to Row 3, the use of 4 splits for adjacency and siblings results in a performance
 that is even worse than using backward adjacency only. This outcome may be attributed to the
 considerable strength of the topological information provided by ASMA, rendering such a complex
 graph structure unnecessary. In Row 4, employing 2 splits with both forward and backward adjacency

486 produces the most favorable result, underscoring the rationale behind our BGIFFN approach. This 487 finding suggests that BGIFFN is well-founded and effective in leveraging topological information. 488

Ablation on BGIFFN design. The gating mechanism in recent neural networks (Ning et al., 2020; 489 Xu et al., 2023) has demonstrated superior performance to the standard feed-forward layer. However, 490 substituting the elementwise add operation in BGIFFN with the Hadamard product results in a 491 significant performance decrease to 0.7076. This may be attributed to the different features of the 492 two branches, as one represents self-position only, and the other aggregates adjacency features. 493 Directly multiplying the two features yields a decrease in performance. Furthermore, we compare 494 our approach with the conventional GCN (Kipf & Welling, 2016) and GAT (Veličković et al., 2018). 495 Both methods lead to a noticeable performance decline, rendering the superiority of our NN-Former. 496

#### 5 CONCLUSION

497

498 499

500

501

506 507

508

517

518

519 520

521

522

527

528

529

We introduce a novel neural architecture representation model. This model unites the strengths of GCN and transformers, demonstrating strong capability in topology modeling and representation learning. We also conclude that different from the intuition on other DAG tasks, sibling nodes 502 significantly affect the extraction of topological information. Our proposed model performs well 503 on accuracy and latency prediction, showcasing model capability and generalization ability. This 504 work may inspire future efforts in neural architecture representation and neural networks on DAG 505 representation.

## REFERENCES

- 509 Mohamed S Abdelfattah, Abhinav Mehrotra, Łukasz Dudziak, and Nicholas D Lane. Zero-cost proxies for lightweight nas. arXiv preprint arXiv:2101.08134, 2021. 510
- 511 Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. Once-for-all: Train one network 512 and specialize it for efficient deployment. arXiv preprint arXiv:1908.09791, 2019. 513
- 514 Yaofo Chen, Yong Guo, Qi Chen, Minli Li, Wei Zeng, Yaowei Wang, and Mingkui Tan. Contrastive 515 neural architecture search with neural architecture comparators. In Proceedings of the IEEE/CVF 516 conference on computer vision and pattern recognition, pp. 9502–9511, 2021.
  - Boyang Deng, Junjie Yan, and Dahua Lin. Peephole: Predicting network performance before training. arXiv preprint arXiv:1712.03351, 2017.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pp. 248–255. Ieee, 2009. 523
- Xiaohan Ding, Xiangyu Zhang, Ningning Ma, Jungong Han, Guiguang Ding, and Jian Sun. Repvgg: 524 Making vgg-style convnets great again. In Proceedings of the IEEE/CVF conference on computer 525 vision and pattern recognition, pp. 13733-13742, 2021. 526
  - Xuanyi Dong and Yi Yang. Nas-bench-201: Extending the scope of reproducible neural architecture search. arXiv preprint arXiv:2001.00326, 2020.
- Zehao Dong, Muhan Zhang, Fuhai Li, and Yixin Chen. Pace: A parallelizable computation encoder 530 for directed acyclic graphs. In International Conference on Machine Learning, pp. 5360–5377. 531 PMLR, 2022. 532
- 533 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas 534 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An 535 image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint 536 arXiv:2010.11929, 2020. 537
- Lukasz Dudziak, Thomas Chau, Mohamed Abdelfattah, Royson Lee, Hyeji Kim, and Nicholas Lane. 538 Brp-nas: Prediction-based nas using gcns. Advances in Neural Information Processing Systems, 33:10480-10490, 2020.

562

568

579

Vijay Prakash Dwivedi and Xavier Bresson. A generalization of transformer networks to graphs.
 *arXiv preprint arXiv:2012.09699*, 2020.

- Mukul Gagrani, Corrado Rainone, Yang Yang, Harris Teague, Wonseok Jeon, Roberto Bondesan,
  Herke van Hoof, Christopher Lott, Weiliang Zeng, and Piero Zappi. Neural topological ordering
  for computation graphs. *Advances in Neural Information Processing Systems*, 35:17327–17339,
  2022.
- Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
   message passing for quantum chemistry. In *International conference on machine learning*, pp. 1263–1272. PMLR, 2017.
- Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017.
- Jianyuan Guo, Kai Han, Han Wu, Yehui Tang, Xinghao Chen, Yunhe Wang, and Chang Xu. Cmt:
   Convolutional neural networks meet vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12175–12185, 2022.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs.
   Advances in neural information processing systems, 30, 2017.
- Kun Jing, Jungang Xu, and Pengfei Li. Graph masked autoencoder enhanced predictor for neural architecture search. In *IJCAI*, pp. 3114–3120, 2022.
- Sam Kaufman, Phitchaya Phothilimthana, Yanqi Zhou, Charith Mendis, Sudip Roy, Amit Sabne, and
   Mike Burrows. A learned performance model for tensor processing units. *Proceedings of Machine Learning and Systems*, 3:387–400, 2021.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks.
   In International Conference on Learning Representations, 2016.
- Bhushan Kotnis, Carolin Lawrence, and Mathias Niepert. Answering complex queries in knowledge
   graphs with bidirectional sequence encoders. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 4968–4977, 2021.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Wei Li, Shaogang Gong, and Xiatian Zhu. Neural graph embedding for neural architecture search. In
   *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 4707–4714, 2020.
- Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 19–34, 2018.
- Liang Liu, Mingzhu Shen, Ruihao Gong, Fengwei Yu, and Hailong Yang. Nnlqp: A multi-platform
   neural network latency query and prediction system with an evolving database. In *Proceedings of the 51st International Conference on Parallel Processing*, pp. 1–14, 2022.
- Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv* preprint arXiv:1608.03983, 2016.
- 586 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint* 587 *arXiv:1711.05101*, 2017.
- Shun Lu, Jixiang Li, Jianchao Tan, Sen Yang, and Ji Liu. Tnasp: A transformer-based nas predictor with a self-evolution framework. *Advances in Neural Information Processing Systems*, 34:15125–15137, 2021.
- Shun Lu, Yu Hu, Peihao Wang, Yan Han, Jianchao Tan, Jixiang Li, Sen Yang, and Ji Liu. Pinat:
   A permutation invariance augmented transformer for nas predictor. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 8957–8965, 2023.

594 595	Renqian Luo, Fei Tian, Tao Qin, Enhong Chen, and Tie-Yan Liu. Neural architecture optimization. <i>Advances in neural information processing systems</i> , 31, 2018.
597 598	Renqian Luo, Xu Tan, Rui Wang, Tao Qin, Enhong Chen, and Tie-Yan Liu. Semi-supervised neural architecture search. <i>Advances in Neural Information Processing Systems</i> , 33:10547–10557, 2020.
599 600 601	Yuankai Luo, Veronika Thost, and Lei Shi. Transformers over directed acyclic graphs. Advances in Neural Information Processing Systems, 36, 2023.
602 603 604	Xuefei Ning, Yin Zheng, Tianchen Zhao, Yu Wang, and Huazhong Yang. A generic graph-based neural architecture encoding scheme for predictor-based nas. In <i>European Conference on Computer Vision</i> , pp. 189–204. Springer, 2020.
605 606 607	Boris T Polyak and Anatoli B Juditsky. Acceleration of stochastic approximation by averaging. <i>SIAM journal on control and optimization</i> , 30(4):838–855, 1992.
608 609	Pranab Kumar Sen. Estimates of the regression coefficient based on kendall's tau. <i>Journal of the American statistical association</i> , 63(324):1379–1389, 1968.
610 611 612 613	Han Shi, Renjie Pi, Hang Xu, Zhenguo Li, James Kwok, and Tong Zhang. Bridging the gap between sample-based and one-shot neural architecture search with bonas. <i>Advances in Neural Information Processing Systems</i> , 33:1808–1819, 2020.
614 615 616	Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Du- mitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 1–9, 2015.
617 618 619	Veronika Thost and Jie Chen. Directed acyclic graph neural networks. <i>arXiv preprint arXiv:2101.07965</i> , 2021.
620 621 622	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. <i>Advances in neural information processing systems</i> , 30, 2017.
623 624 625 626	Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In <i>International Conference on Learning Representations</i> , 2018.
627 628 629	Wei Wen, Hanxiao Liu, Yiran Chen, Hai Li, Gabriel Bender, and Pieter-Jan Kindermans. Neural pre- dictor for neural architecture search. In <i>Computer Vision–ECCV 2020: 16th European Conference</i> , <i>Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIX</i> , pp. 660–676. Springer, 2020.
630 631 632	Colin White, Willie Neiswanger, and Yash Savani. Bananas: Bayesian optimization with neural architectures for neural architecture search. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pp. 10293–10301, 2021.
634 635	Felix Wu, Angela Fan, Alexei Baevski, Yann N Dauphin, and Michael Auli. Pay less attention with lightweight and dynamic convolutions. <i>arXiv preprint arXiv:1901.10430</i> , 2019.
636 637 638 639	Zhanghao Wu, Paras Jain, Matthew Wright, Azalia Mirhoseini, Joseph E Gonzalez, and Ion Stoica. Representing long-range context for graph neural networks with global attention. <i>Advances in</i> <i>Neural Information Processing Systems</i> , 34:13266–13279, 2021.
640 641	Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In <i>International Conference on Learning Representations</i> , 2018.
642 643 644 645	Ruihan Xu, Haokui Zhang, Wenze Hu, Shiliang Zhang, and Xiaoyu Wang. Parcnetv2: Oversized kernel with enhanced attention. In <i>Proceedings of the IEEE/CVF International Conference on Computer Vision</i> , pp. 5752–5762, 2023.
646 647	Yixing Xu, Yunhe Wang, Kai Han, Yehui Tang, Shangling Jui, Chunjing Xu, and Chang Xu. Renas: Relativistic evaluation of neural architecture search. In <i>Proceedings of the IEEE/CVF conference</i> on computer vision and pattern recognition, pp. 4411–4420, 2021.

648	Shen Yan Yu Zheng Wei Ao Xiao Zeng and Mi Zhang Does unsupervised architecture representa-
649	tion learning help neural architecture search? Advances in Neural Information Processing Systems.
650	33:12486–12498, 2020.
651	

- Yun Yi, Haokui Zhang, Wenze Hu, Nannan Wang, and Xiaoyu Wang. Nar-former: Neural architecture representation learning towards holistic attributes prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7715–7724, 2023.
- Yun Yi, Haokui Zhang, Rong Xiao, Nannan Wang, and Xiaoyu Wang. Nar-former v2: Rethinking
   transformer for universal neural network representation learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- 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 *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021. URL https://openreview. net/forum?id=OeWooOxFwDa.
- Chris Ying, Aaron Klein, Eric Christiansen, Esteban Real, Kevin Murphy, and Frank Hutter. Nas bench-101: Towards reproducible neural architecture search. In *International conference on machine learning*, pp. 7105–7114. PMLR, 2019.
- Jiaxuan You, Jure Leskovec, Kaiming He, and Saining Xie. Graph structure of neural networks. In International Conference on Machine Learning, pp. 10881–10891. PMLR, 2020.
- Chris Zhang, Mengye Ren, and Raquel Urtasun. Graph hypernetworks for neural architecture search.
   *arXiv preprint arXiv:1810.05749*, 2018.

Li Lyna Zhang, Shihao Han, Jianyu Wei, Ningxin Zheng, Ting Cao, Yuqing Yang, and Yunxin Liu. Nn-meter: Towards accurate latency prediction of deep-learning model inference on diverse edge devices. In *Proceedings of the 19th Annual International Conference on Mobile Systems, Applications, and Services*, pp. 81–93, 2021.

# 702 A METHODS DETAILS

#### 704 A.1 IMPLEMENTATION FOR ASMA 705

We present Python-style code for calculating the attention matrix in the ASMA module in Listing 1. 706 ASMA is motivated by the importance of sibling nodes. In the accuracy prediction, sibling nodes 707 provide complementary features, such as parallel 1x1 and 3x3 convolutions extracting pixel features 708 and local aggregations, respectively. Although the two nodes are neither connected nor reachable through transitive closure, their information can influence each other. This conclusion has been 710 studied in works such as Inception (Szegedy et al., 2015) and RepVGG (Ding et al., 2021). In latency 711 prediction, sibling nodes can run in parallel. For example, if two parallel 1x1 convolutions are merged 712 into one, it takes only one CUDA kernel and fully utilizes parallel computing. Hence it is reasonable 713 for ASMA to fuse the sibling nodes information directly. 714

Listing 1: Calculating the attention matrix in ASMA.

```
def attention_matrix(Q, K, A):
    # Q: query, K: key, A: adjacency matrix
    # Calculate the attention scores
    attn = torch.matmul(Q, K.mT) / math.sqrt(Q.size(-1))
    # Prepare attention masks
    pe = torch.stack([A, A.mT, A.mT @ A, A @ A.mT], dim=1)
    pe = pe + torch.eye(L, dtype=A.dtype, device=A.device)
    # Apply masking
    attn = attn.masked_fill(pe == 0, -torch.inf)
    # Softmax operation
    attn = F.softmax(attn, dim=-1)
    return attn
```

729

730 731

732

739

740

715

716

717

718

719

720

721

722

723

724

To implement the masking operation, the values at the non-zero positions remain unchanged, while the other values are set to minus infinity. Consequently, the softmax operation on these masked values results in zeroes.

#### A.2 PROOF OF SIBLING NODES IDENTIFICATION

In the paper we use  $A^T A$  to represent sibling nodes that share a same successor. Here we provide a trivial proof.  $A_{ij} = 1$  denotes there is a directed edge linked from node *i* to node *k*. Thus  $A_{ki}^T = 1$ denotes that there is a directed edge linked from node *i* to node *k*. Thus  $(A^T A)_{kj} = \sum_v A_{kv}^T A_{vj} \ge$  $A_{ki}^T A_{ij} = 1$ , which denotes that node *k* and node *j* share a same successor *i*. Similar to  $AA^T$ , where  $(AA^T)_{kj} \ge 1$  if node *k* and node *j* share a same predecossor.

#### A.3 PROOF OF BI-DIRECTIONAL GRAPH ISOMORPHISM FEED-FORWARD NETWORK

We begin by summarizing the BIGFFN as the common form of message-passing GNNs, and then prove the isomorphism property. Modern message-passing GNNs follow a neighborhood aggregation strategy, where we iteratively update the representation of a node by aggregating representations of its neighbors. To make comparison with modern GNNs, we follow the same notations, where the feature of node v is denoted as  $h_v$ . The *l*-th layer of a GNN is composed of aggregation and combination operation:

$$a_v^{(l)} = \operatorname{AGGREGATE}\left(h_u^{(l)} : u \in \mathcal{N}(v)\right), h_v^{(l)} = \operatorname{COMBINE}\left(h_v^{(l-1)}, a_v^{(l)}\right),$$
(14)

where  $h_v^{(l)}$  is the feature vector of node v at the *l*-th iteration/layer. In our cases, the graph is directional, thus the neighborhood  $\mathcal{N}(v)$  is also divided into forward propagation nodes  $\mathcal{N}^+(v)$  and backward propagation nodes  $\mathcal{N}^-(v)$ :

747 748

$$a_v^{(l)} = \operatorname{AGGREGATE}\left(h_u^{(l-1)} : u \in \mathcal{N}^+(v) \cup \mathcal{N}^-(v)\right).$$
(15)

The AGGREGATE function in BGIFFN is defined as a matrix multiplication followed by concatenation:

AGGREGATE : 
$$H \mapsto \text{Concat} \left( AHW^+, A^T HW^- \right)$$
, (16)

where  $W^+$  and  $W^-$  are the linear transform for forward and backward propagation, respectively. This is equivalent to a bidirectional neighborhood aggregation followed by a concatenation operation:

762

763 764 765

769

770

771

773 774 775

776

777 778

784

802 803

808 809

$$a_v^{(l)} = \operatorname{Concat}\left(\sum_{u \in \mathcal{N}^+(v)} h_u^{(l-1)} \boldsymbol{W}^+, \sum_{u \in \mathcal{N}^-(v)} h_u^{(l-1)} \boldsymbol{W}^-\right),\tag{17}$$

and the COMBINE function is defined as follows:

$$h_v^{(l)} = \operatorname{ReLU}\left(h_v^{(l-1)}\boldsymbol{W}_1 + a_v^{(l)}\right)\boldsymbol{W}_2.$$
(18)

We quote Theorem 3 in (Xu et al., 2018). For simple reference, we provide the theorem in the following:

**Theorem 1** (Theorem 3 in (Xu et al., 2018)). With a sufficient number of GNN layers, a GNN  $\mathcal{M} : \mathcal{G} \mapsto \mathbb{R}^d$  maps any graphs  $G_1$  and  $G_2$  that the Weisfeiler-Lehman test of isomorphism decides as non-isomorphic, to different embeddings if the following conditions hold:

a)  $\mathcal{T}$  aggregates and updates node features iteratively with

$$h_{v}^{(l)} = \phi\left(h_{v}^{(l-1)}, f\left(\left\{h_{v}^{(l-1)} : u \in \mathcal{N}(v)\right\}\right)\right),$$
(19)

where the function f, which operates on multisets, and  $\varphi$  are injective.

b)  $\mathcal{T}$ 's graph-level readout, which operates on the multiset of node features  $\left\{h_v^{(l)}\right\}$ , is injective.

Please refer to (Xu et al., 2018) for the proof. In our cases, the difference lies in condition a), where our tasks use directed acyclic graphs. Thus we modify condition a) as follows:

**Theorem 2** (Modified condition for undirected graph). *a*)  $\mathcal{T}$  aggregates and updates node features *iteratively with* 

$$h_{v}^{(l)} = \phi\left(h_{v}^{(l-1)}, f\left(\left\{h_{v}^{(l-1)} : u \in \mathcal{N}^{+}(v) \cup \mathcal{N}^{-}(v)\right\}\right)\right),$$
(20)

where the function f, which operate on multisets, and  $\varphi$  are injective.

The proof is trivial, as it turns back to the original undirected graph. Following the Corollary 6 in (Xu et al., 2018), we can build our bidirectional graph isomorphism feed-forward network:

**Corollary 1** (Corollary 6 in (Xu et al., 2018)). Assume  $\mathcal{X}$  is countable. There exists a function  $f: \mathcal{X} \to \mathbb{R}^n$  so that for infinitely many choices of  $\epsilon$ , including all irrational numbers,  $h(c, X) = (1 + \epsilon) \cdot f(c) + \sum_{x \in X} f(x)$  is unique for each pair (c, X), where  $c \in \mathcal{X}$  and  $X \subset \mathcal{X}$  is a multiset of bounded size. Moreover, any function g over such pairs can be decomposed as  $g(c, X) = \varphi((1 + \epsilon) \cdot f(c) + \sum_{x \in X} f(x))$  for some function  $\varphi$ .

<sup>795</sup> In our cases,  $\epsilon$  is substituted by a linear transform with weights  $W_1$ . f is the aggregation function, <sup>796</sup> and  $\varphi$  is the combine function. There exist choices of f and  $\varphi$  that are injective, thus the conditions <sup>797</sup> are satisfied.

Furthermore, our BGIFFN distinguishes the forward and backward propagation, yielding stronger capability in modeling graph topology. Our method corresponds to a stronger "directed WL test", which applies a predetermined injective function z to update the WL node labels  $k_v^{(l)}$ :

$$k_{v}^{(l)} = z\left(k_{v}^{(l)}, \left\{k_{v}^{(l)} : u \in \mathcal{N}^{+}(v)\right\}, \left\{k_{v}^{(l)} : u \in \mathcal{N}^{-}(v)\right\}\right),$$
(21)

and the condition is modified as:

**Theorem 3** (Modified condition for directed graph). *a*) T aggregates and updates node features *iteratively with* 

$$h_{v}^{(l)} = \phi\left(h_{v}^{(l-1)}, f\left(\left\{h_{v}^{(l-1)} : u \in \mathcal{N}^{+}(v)\right\}\right), g\left(\left\{h_{v}^{(l-1)} : u \in \mathcal{N}^{-}(v)\right\}\right)\right),$$
(22)

where the function f and g, which operate on multisets, and  $\varphi$  are injective.

810 *Proof.* The proof is a trivial extension to Theorem 1. Let  $\mathcal{T}$  be a GNN where the condition holds. 811 Let  $G_1$  and  $G_2$  be any graphs that the directed WL-test (which means propagating on the directed 812 graph) decides as non-isomorphic at iteration L. Because the graph-level readout function is injective, 813 it suffices to show that  $\mathcal{T}$ 's neighborhood aggregation process embeds  $G_1$  and  $G_2$  into different 814 multisets of node features with sufficient iterations. We will show that for any iteration l, there always exists an injective function  $\varphi$  such that  $h_v^{(k)} = \varphi\left(k_v^{(l)}\right)$ . This holds for l = 0 because the initial node 815 816 features are the same for WL and GNN  $k_v^{(0)} = h_v^{(0)}$ . So  $\varphi$  could be the identity function for k = 0. 817 Suppose this holds for iteration k - 1, we show that it also holds for l. Substituting  $h_v^{(l-1)}$  with 818  $\varphi\left(h_v^{(l-1)}\right)$  gives us: 819 820

$$h_{v}^{(l)} = \phi\left(\varphi\left(h_{v}^{(l-1)}\right), f\left(\left\{\varphi\left(h_{v}^{(l-1)}\right) : u \in \mathcal{N}^{+}(v)\right\}\right), g\left(\left\{\varphi\left(h_{v}^{(l-1)}\right) : u \in \mathcal{N}^{-}(v)\right\}\right)\right),$$
(23)

Since the composition of injective functions is injective, there exists some injective function  $\psi$  so that

$$h_{v}^{(l)} = \psi\left(h_{v}^{(l-1)}, \left\{h_{v}^{(l-1)} : u \in \mathcal{N}^{+}(v)\right\}, \left\{h_{v}^{(l-1)} : u \in \mathcal{N}^{-}(v)\right\}\right).$$
(24)

Then we have

821 822 823

824 825 826

827 828 829

837

838 839

840

841

842

843 844

845

$$h_{v}^{(l)} = \psi \circ z^{-1} z \left( k_{v}^{(l)}, \left\{ k_{v}^{(l)} : u \in \mathcal{N}^{+}(v) \right\}, \left\{ k_{v}^{(l)} : u \in \mathcal{N}^{-}(v) \right\} \right),$$
(25)

830 and thus  $\varphi = \psi \circ z^{-1}$  is injective because the composition of injective functions is injective. Hence 831 for any iteration l, there always exists an injective function  $\varphi$  such that  $h_v^{(l)} = \varphi(h_v^{(l-1)})$ . At the 832 L-th iteration, the WL test decides that  $G_1$  and  $G_2$  are non-isomorphic, that is the multisets  $k_n^L$  are 833 different for  $G_1$  and  $G_2$ . The graph neural network  $\mathcal{T}$ 's node embeddings  $\left\{h_v^{(L)}\right\} = \left\{\varphi\left(k_v^{(L)}\right)\right\}$ 834 835 must also be different for  $G_1$  and  $G_2$  because of the injectivity of  $\varphi$ . 836

#### A.4 IMPLEMENTATION FOR BGIFFN

# x: node features, A: adjacency matrix

combine = F.relu(x @ W\_1 + aggregate) @ W\_2

We present Python-style code for the BGIFFN module in Listing 2. BGIFFN is intended to extend Graph Isomorpsim to the bidirectional modeling of DAGs. It extracts the topological features simply and effectively, assisting the Transformer backbone in learning the DAG structure. Various works use convolution to enhance FFN in vision (Guo et al., 2022) and language tasks (Wu et al., 2019). It is reasonable for BGIFFN to assist Transformer in neural predictors.

```
Listing 2: Calculation for BGIFFN.
```

# W\_1, W\_forward, W\_backward, W\_2: the weight for linear transform

aggregate = torch.cat((A @ x @ W\_forward, A.mT @ x @ W\_backward),

```
def bgiffn(x, A, W_1, W_forward, W_backward, W_2):
846
847
848
849
850
```

855 856

858

859

861

#### **EXPERIMENT DETAILS** В

dim=-1)

return combine

We present implementation details of our proposed NN-Former. For accuracy prediction, we show the experiment settings on NAS-Bench-101 in Section B.1.1 and NAS-Bench-201 in Section B.1.2. For latency prediction, we show the experiment settings on NNLQ (Liu et al., 2022) in Section B.2.1.

**B.1** ACCURACY PREDICTION

For the network input, each operation type is represented by a 32-dimensional vector using one-hot 862 encoding. Subsequently, this encoding is converted into a 160-channel feature by a linear transform 863 and a layer normalization. The model contains 12 transformer blocks commonly employed in vision

864 transformers (Dosovitskiy et al., 2020). Each block comprises ASMA and BGIFFN modules. The 865 BGIFFN has an expansion ratio of 4, mirroring that of a vision transformer. The output class token is 866 transformed into the final prediction value through a linear layer. Initialization of the model follows 867 a truncated normal distribution with a standard deviation of 0.02. During training, Mean Squared 868 Error (MSE) loss is utilized, alongside other augmentation losses as outlined in NAR-Former (Yi et al., 2023) with  $\lambda_1 = 0.2$  and  $\lambda_2 = 1.0$ . The model is trained for 3000 epochs in total. A warmup (Goyal et al., 2017) learning rate from 1e-6 to 1e-4 is applied for the initial 300 epochs, and cosine 870 annealing (Loshchilov & Hutter, 2016) is adopted for the remaining duration. AdamW (Loshchilov 871 & Hutter, 2017) with a coefficient (0.9, 0.999) is utilized as the optimizer. The weight decay is set 872 to 0.01 for all the layers except that the layer normalizations and biases use no weight decay. The 873 dropout rate is set to 0.1. We use the Exponential Moving Average (EMA) (Polyak & Juditsky, 1992) 874 with a decay rate of 0.99 to alleviate overfitting. Each experiment takes about 1 hour to train on an 875 RTX 3090 GPU. 876

877 878

### B.1.1 EXPERIMENTS ON NAS-BENCH-101.

NAS-Bench-101 (Ying et al., 2019) provides the performance of each architecture on CIFAR10 (Krizhevsky et al., 2009). It is an operation-on-node (OON) search space, which means nodes
represent operations, while edges illustrate the connections between these nodes. Following the
approach of TNASP (Lu et al., 2021), we utilize the validation accuracy from a single run as the
target during training, and the mean test accuracy over three runs is used as ground truth to assess the
Kendall's Tau (Sen, 1968). The metrics on the test set are computed using the final epoch model, the
top-performing model, and the best Exponential Moving Average (EMA) model on the validation set.
The highest-performing model is documented.

886 887

#### B.1.2 EXPERIMENTS ON NAS-BENCH-201.

NAS-Bench-201 offers three sets of results for each architecture, corresponding to CIFAR-10, CIFAR 100, and ImageNet-16-120. This study focuses on the CIFAR-10 dataset, consistent with the setup in
 TNASP (Lu et al., 2021).

892 NAS-Bench-201 (Dong & Yang, 2020) is originally operation-on-edge (OOE) search space, while 893 we transformed the dataset into the OON format. NAS-Bench-201 contains the performance of 894 each architecture on three datasets: CIFAR-10 (Krizhevsky et al., 2009), CIFAR-100 (Krizhevsky 895 et al., 2009), and ImageNet-16-120 (a downsampled subset of ImageNet (Deng et al., 2009)). We 896 use the results on CIFAR-10 in our experiments following previous TNASP (Lu et al., 2021), NAR-Former (Yi et al., 2023) and PINAT (Lu et al., 2023). In the preprocessing, we drop the useless 897 898 operations taht only have zeroized input or output. The metrics on the test set are computed using the final epoch model, the top-performing model, and the best Exponential Moving Average (EMA) 899 model on the validation set. The highest-performing model is documented. 900

901 As for the results in the 10% setting, we argue that these results are not a good measurement. 902 Concretely, the predictors are trained on the validation accuracy of NAS-Bench-201 networks, and 903 evaluated on the test accuracy. We calculate Kendall's Tau between ground truth validation accuracy and test accuracy on this dataset which is 0.889. It indicates an unneglectable gap between the 904 predictors' training and testing. Thus the results around and higher than 0.889 are less valuable to 905 reflect the performance of predictors. For further studies, we also provide a new setting for this 906 dataset. Both training and evaluation are conducted on the test accuracy of NAS-Bench-201 networks, 907 and the training samples are dropped during evaluation. This setting has no gap between the training 908 and testing distribution. As shown in Table 8, our methods surpass both NAR-Former (Yi et al., 2023) 909 and NAR-Former V2 (Yi et al., 2024), showcasing the strong capability of our NN-Former. 910

- 910
- 912 B.2 LATENCY PREDICTION
- 913 B.2.1 EXPERIMENTS ON NNLQ. 914

There are two scenarios on latecny prediction on NNLQ (Liu et al., 2022). In the first scenario, the training set is composed of the first 1800 samples from each of the ten network types, and the remaining 200 samples for each type are used as the testing set. The second scenario comprises ten sets of experiments, where each set uses one type of network as the test set and the remaining Table 8: Accuracy prediction results on NAS-Bench-201 (Dong & Yang, 2020) when the training
and testing data follow the same distribution. We use different proportions of data as the training
set and report Kendall's Tau on the whole dataset.

Method	Publication		Training Samples 10% (1563)
NAR-Former (Yi et al., 2023)	CVPR 2023		0.910
NAR-Former V2 (Yi et al., 2024) NN-Former (Ours)	NeurIPS 202	.3	0.921 <b>0.935</b>

929 nine types serve as the training set. The network input is encoded in a similar way as NAR-Former 930 V2 (Yi et al., 2024). Each operation is represented by a 192-dimensional vector, with 32 dimensions 931 of one-hot operation type encoding, 80 dimensions of sinusoidal operation attributes encoding, 932 and 80 dimensions of sinusoidal feature shape encoding. Subsequently, this encoding is converted into a 512-channel feature by a linear transform and a layer normalization. The model contains 2 933 transformer blocks, the same as NAR-Former V2 (Yi et al., 2024). Each block comprises ASMA 934 and BGIFFN modules. The BGIFFN has an expansion ratio of 4, mirroring that of a common 935 transformer (Dosovitskiy et al., 2020). The output features are summed up and transformed into the 936 final prediction value through a 2-layer feed-forward network. Initialization of the model follows 937 a truncated normal distribution with a standard deviation of 0.02. During training, Mean Squared 938 Error (MSE) loss is utilized. The model is trained for 50 epochs in total. A warm-up (Goyal et al., 939 2017) learning rate from 1e-6 to 1e-4 is applied for the initial 5 epochs, and a linear decay scheduler 940 is adopted for the remaining duration. AdamW (Loshchilov & Hutter, 2017) with a coefficient (0.9, 941 0.999) is utilized as the optimizer. The weight decay is set to 0.01 for all the layers except that the 942 layer normalizations and biases use no weight decay. The dropout rate is set to 0.05. We also use 943 static features as NAR-Former V2 (Yi et al., 2024). Each experiment takes about 4 hours to train on 944 an RTX 3090 GPU.

945 946

947

948 949

950 951

952

953

954

955 956

957

958

# C EXTENSIVE EXPERIMENTS

# C.1 ABLATION ON HYPERPARAMETERS

This work adopts a Transformer as the backbone, and the hyperparameters of Transformers have been well-settled in previous research. This article follows the common training settings (from NAR-Former) and has achieved good results. Apart from these hyperparameters, we provide an ablation on the number of channels and layers in the predictor as shown in Table 9:

Table 9: Ablation studies on hyperparameters. All experiments are conducted on the NAS-Bench-101 (Ying et al., 2019) with the 0.04% training set. (a) Ablation study on the number of channels. (b) Ablation study on the number of transformer layers.

(a)		(b)	(b)			
Num of Channe	ls   KT	Num of Layers	KT			
64	0.748	6	0.744			
128	0.758	9	0.760			
160 (Ours)	0.765	12 (Ours)	0.765			

966 967 968

969 970

#### C.2 COMPARISON WITH ZERO-COST PREDICTORS

971 Zero-cost proxies are lightweight NAS methods, but they performs not as well as the model-based neural predictors.

NAS search space	NAS-Bench-101↑	NAS-Bench-201 <sup>†</sup>
grad_norm (Abdelfattah et al., 2021)	0.20	0.58
snip (Abdelfattah et al., 2021)	0.16	0.58
NN-Former	0.71	0.80

Table 10: Comparison with zero-cost predictors.

#### MODEL COMPLEXITY D

#### 982 D.1 THEORETICAL ANALYSIS 983

Our ASMA method has less or equal computational complexity than the vanilla attention. On the dense graph, the vanilla self-attention has a complexity of  $O(N^2)$  where N denotes the number of nodes. With the sibling connection preprocessed, our ASMA also has a complexity of  $O(N^2)$ . On sparse graphs, the vanilla self-attention is still a global operation thus the complexity is also  $O(N^2)$ . Our ASMA has a complexity of O(NK), where K is the average degree and  $K \ll N$  on sparse graphs. In practical applications, sparse graphs are common thus our method is efficient. The latency prediction experiments in the paper show that our predictor can cover the DAGs from 20 200 nodes, which is applicable for practical use.

991 992 993

994

#### D.2 INFERENCE SPEED

995 We report the parameters and the latency, memory, and training time on a single RTX 3090 in Table 11. Our method has comparable inference latency, memory usage, and training time compared 996 to NAR-Former Yi et al. (2023), indicating that the improvement brought by our method is solid. 997

998 Compared to the tremendous time spent training candidate architectures, the time of training neural 999 predictors is neglectable. Therefore, the computational resources consumed by predictors do not 1000 affect practical applications. In the experiments on NNLQ, our method can make predictions for networks with from 20 to 200 layers, which encompasses the size of most practical models. It 1001 indicates that our method can be applied to practical use. 1002

Table 11: Caption	on
-------------------	----

1006	Predictor	Params (M)	Latency (ms)	Memory (GB)	Training Time (h)
007	NAR-Former	4.8	10.31	0.58	0.7
800	NN-Former w/o ASMA	4.9	11.21	0.67	0.8
009	NN-Former w/o BGIFFN	3.7	10.17	0.60	0.7
010	NN-Former	4.9	11.53	0.67	0.8

1011 1012 1013

1014

1003 1004

#### SIBLING NODES MODELING ON GENERAL DAG TASKS E

1015 Our method is dedicated to neural architecture representation learning. Experiments have shown that 1016 the sibling nodes provide valuable insights into accuracy prediction and latency prediction. However, 1017 considering the nature of our method is DAG modeling, one question comes out: is it possible for the 1018 sibling property to generalize to other DAG tasks? Research on other DAG tasks is beyond the scope 1019 of this study. However, we can provide theoretical analysis with validation experiments. The DAG 1020 tasks are divided into two groups by the importance of sibling nodes:

1021 One is that sibling nodes are important. Citation prediction is a situation where sibling nodes play 1022 a crucial role. Two papers that cite the same paper might follow similar motivations, methods, or 1023 experiments. Similarly, two papers cited by the same paper may also have these in common. 1024

The other is that siblings are not as important. For example, Abstract Syntax Tree (AST) uses 1025 syntactic construct to aggregate the successors, while siblings do not make practical sense.

972

981

984

985

986

987

988

989

990

1026 1027	Table 12: Sibli	Table 12: Sibling perspective on other DAG tasks.					
1028	Model	Cora(%)↑	ogbg-code2(%)↑				
1029	DAC Transformer						
1030	DAG Transformer	sibling 87.39	19.0				
1031	DAG Hallstoffler +	sibiling   00.14	10.9				
1032							
1033	We also experimented on how sibling	s affect the results in	these tasks We use	e the DAG Trans			
1034	former (Luo et al. 2023) as the baseling	and add a sibling atte	ention mask without	careful calibration			
1035	The results in the Table 12 show that s	sibling nodes play a c	rucial role in Cora n	ode classification			
1036	(Citation prediction), while they are not	t as important in ogbg	-code2 (AST predict	ion).			
1037		1 00					
1038							
1039							
1040							
1041							
1042							
1043							
1044							
1045							
1046							
1047							
1048							
1049							
1050							
1051							
1052							
1053							
1055							
1055							
1057							
1058							
1059							
1060							
1061							
1062							
1063							
1064							
1065							
1066							
1067							
1068							
1069							
1070							
1071							
1072							
1073							
1074							
1075							
1075							
1077							
1079							

# Table 12: Sibling perspective on other DAG tasks