

---

# Understanding and Simplifying Architecture Search in Spatio-Temporal Graph Neural Networks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1        Compiling together spatial and temporal modules via a unified framework, Spatio-  
2        Temporal Graph Neural Networks (STGNNs) have been popularly used in the  
3        multivariate spatio-temporal forecasting task, e.g. traffic prediction. After the  
4        numerous propositions of manually designed architectures, researchers show in-  
5        terest in the Neural Architecture Search (NAS) of STGNNs. Existing methods  
6        suffer from two issues: (1) hyperparameters like learning rate, channel size cannot  
7        be integrated into the NAS framework, which makes the model evaluation less  
8        accurate, potentially misleading the architecture search (2) the current search space,  
9        which basically mimics Darts-like methods, is too large for the search algorithm  
10       to find a sufficiently good candidate. In this work, we deal with both issues at the  
11       same time. We first re-examine the importance and transferability of the training  
12       hyperparameters to ensure a fair and fast comparison. Next, we set up a framework  
13       that disentangles architecture design into three disjoint angles according to how  
14       spatio-temporal representations flow and transform in architectures, which allows  
15       us to understand the behavior of architectures from a distributional perspective.  
16       This way, we can obtain good guidelines to reduce the STGNN search space and  
17       find state-of-the-art architectures by simple random search. As an illustrative exam-  
18       ple, we combine these principles with random search which already significantly  
19       outperforms both state-of-the-art hand-designed models and recently automatically  
20       searched ones.<sup>1</sup>

## 21    1 Introduction

22    Multivariate forecasting is a crucial task commonly encountered in our daily life. In addition to  
23    classic time series forecasting, spatial correlation has been recently leveraged for more accurate  
24    prediction, titled Spatio-Temporal Forecasting. A typical and important example is traffic prediction  
25    which forecasts future traffic status (e.g., volume and speed) based on history data (usually the  
26    same as prediction targets), moving one step further towards smart transportation and intelligent  
27    city [40, 26, 24]. Using learning models to better exploit spatio-temporal information is the key  
28    challenge [15, 18].

29    Early works use classic time series forecasting techniques like ARIMA [39] and LSTM [11] but  
30    do not explicitly consider spatial relation. Motivated by the power of Graph Neural Networks  
31    (GNNs) [27, 16] modeling relational data, two representative works [37, 20] firstly combine GNN  
32    and temporal modeling techniques and propose Spatio-Temporal Graph Neural Networks (STGNN).  
33    STGNNs simultaneously extract spatial correlation by GNN module and temporal correlation by,  
34    like Convolutional Neural Network (CNN) [17] or Gated Recurrent Unit (GRU) [6], and they have

---

<sup>1</sup>Our code is available at <https://github.com/AutoML-Research/SimpleSTG>.

shown promising performance in traffic prediction. Subsequently, many variants are proposed to improve upon above pioneer works from various angles. Examples are STGCN [37], DCRNN [20], MTGNN [33], AGCRN [1], and STFGCN [19].

As STGNNs go more complex, designing better architectures by hand can go beyond human expertise because it remains difficult for researchers to understand which architecture configurations work better. More recently, Neural Architecture Search(NAS) methods have gained much attention due to their effectiveness in various research directions [31, 9, 22]. NAS researchers leverage the expertise in certain domains to propose efficient search space and strategy. Similarly, two STGNN NAS methods have been introduced as well AutoCTS [32] and AutoSTG [25], which basically follow Darts-like differentiable search space [22] with a few customized operations to adapt on STGNN models. However, since there is a lack of *systematic architecture understanding* in the STGNN community, the STGNN NAS methods still suffer from a huge search space, hindering their efficiency and effectiveness. Moreover, all the STGNN NAS methods do not consider the impact of *hyperparameters*, leading to probable evaluation bias.

Above works usually propose a novel architecture (manually or by search) and study this single model in an isolated way. Motivated by recent works that study a series of related works in hindsight to produce general principles [7, 12, 3, 36, 41], we aim at obtaining useful understandings to bring forth the next generation of STGNNs, which is an important missing piece in the community.

However, systematic understanding of STGNN architectures is not an easy task. First, the evaluation of architectures is expensive as hyperparameter tuning is required. Training hyperparameters in existing works are often tuned with a small grid inherited from prior studies. But different architectures may not prefer the same setting of hyperparameters and a small grid search cannot fully explore the hyperparameters space [36, 3]. Currently, there is no existing methods that can evaluate and compare architectures in a cheap and fair way. Second, the architecture space is huge and complicated. Classically, control variable methods, which analyze improvements with and without certain microscopic architecture designs, are popularly used [1, 28]. Reusing such methods here could lead to biased interpretations as they fail to explore different architectures from a distributional perspective.

In this work, we take a step back and revisit the STGNNs. Inspired by the representative literature as illustrated in Figure 1, we propose a disentangled framework composed of disjoint factors regarding the architecture designs. This framework includes design choices of temporal/spatial modules, spatio-temporal order and skip connection. Through the lens of this framework, we dig in depth the principles of different architectural choices. We summarize and quantify empirically the influence of each part. Based on our understanding, we could easily find new STGNN models with a simple method adapted from random search.

Our main contributions are as follows:

- We propose a disentangled framework of Spatio-Temporal Graph Neural Networks containing the designs of spatial/temporal modules, the spatio-temporal order and the skip connection.
- To allow fair and efficient evaluation of different model configurations, we study a much larger hyperparameter setting and quantify the importance of each choice, which helps to reduce significantly the hyperparameter space.
- Through comprehensive experiments and distributional analysis, we conclude fundamental design principles of the skip connection choices and the spatio-temporal order, quantified by mathematically motivated measures. The design principles form naturally cherry regions, where architectures are more likely to perform well.
- Based on our in depth understanding of architectural principles and training hyperparameters, we propose a simple method to obtain new STGNN models more easily and effectively.

## 2 Problem Definition

A graph representing the topology network is denoted as  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}, A\}$ , where  $\mathcal{V}$  is the set of  $|\mathcal{V}| = N$  nodes,  $\mathcal{E}$  is the edge set and  $A$  is the adjacency matrix of shape  $\mathbb{R}^{N \times N}$ . We have then a multivariate feature tensor of nodes  $X \in \mathbb{R}^{T \times N \times D}$  representing the temporal graph signals.  $T$  is the

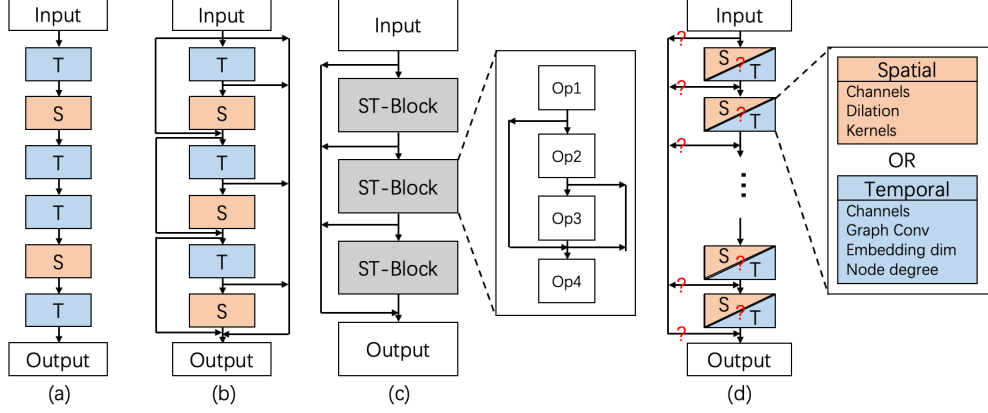


Figure 1: Architecture plot of representative literatures. (a) STGCN (b) MTGNN (c) AutoCTS (d) Our framework. Red blocks represent spatial modules (S). Blue blocks represent temporal modules (T). For AutoCTS, each block composes of four operations and each operation could be spatial or temporal module.

total timestamps.  $D$  is the number of features per node at each timestamp. Denote also  $X_t \in \mathbb{R}^{N \times D}$  the features at one timestamp  $t$  and  $X_{t:(t+T')} \in \mathbb{R}^{T' \times N \times D}$  the features of  $T'$  timestamps.

In spatio-temporal traffic prediction task, the goal is to predict  $Q$  future steps based on  $P$  past steps observations of traffic features [18]. Due to the popularity and superior performance, we use STGNNs as the prediction model. Following the mentioned literature [33, 5, 32], we consider  $P = Q = 12$ .

In order to understand the generalizable principles of STGNN architectures in a **tractable** way, we propose the following **distributional** bi-level formulation:

$$\min_{\theta} \mathbb{E}_{\alpha \sim \mathcal{P}_{\theta}(\alpha)} [\mathcal{L}(w^*(\alpha), \alpha; D_{\text{val}})] \quad \text{s.t.} \quad w^*(\alpha) = \arg \min_w \mathcal{L}(w, \alpha; D_{\text{train}}),$$

where  $\alpha$  is the model architecture, encoded in any raw format;  $w$  is the trainable parameters depending on  $\alpha$ ;  $\mathcal{P}$  is the architectural distribution parameterized by  $\theta$ ;  $\mathcal{L}$  is the evaluation metric such as mean squared error;  $D_{\text{val}}$ ,  $D_{\text{train}}$  are validation and training data. The usage of  $D_{\text{val}}$  in the outer optimization and  $D_{\text{train}}$  in the inner optimization is commonly adopted in the NAS community in order to separate conceptually two levels of problem [22, 4]. In our work, we follow this paradigm.

The distributional nature lies in the minimization of expected loss, similar to [4, 34]. And since the raw architecture encoding is not compact and informative enough, we turn this task into a tractable one by proposing parameterized measures  $\theta$  to map from raw architectures  $\alpha$  to their measures  $\theta$  and finally to the distributional predictive performances.

Generally speaking, the architecture  $\alpha$  can capture certain prior knowledge on the traffic prediction task. For example, the order of how spatial/temporal representations are processed, how different levels of spatio-temporal representations could be fused. All these prior knowledge should be data-dependent and can largely influence the architecture choice.

Thus, understanding of architectures leads to better exploit the prior knowledge and subsequently obtain better prediction performance. As discussed in Section 1, existing works cannot be used since the architecture’s space is large and complex and a fair evaluation is expensive. In the sequel, we propose a disentangled framework and choose appropriate measures to evaluate architectures on observations generated from the framework.

### 3 A Disentangled Framework

To motivate such a framework, we take a step back and revisit representative works’ architectures in Figure 1(a)-(d). We remove unnecessary microscopic details and highlight their overall architectures. STGCN and AGCRN are simplified to compile 6 blocks of Spatial/Temporal modules in different orders (resp. T-S-T-T-S-T and S-T-S-T-S-T). MTGNN compiles 6 blocks in another order (T-S-T-S-T-S) and in addition, adds skip connections. AutoCTS consists of multiple ST-Blocks, each of

Table 1: Summarization of the proposed architecture framework

Type	Name	Detail
Spatio-temporal order	Arrangement of spatial/temporal modules in arbitrary order	
Skip connection	Any jumping between two modules including input and output	
Spatial module	Channels	16,32,64
	Graph convolution	Kipf GCN, Cheb GCN, Mixhop GCN
	Embedding dim	10,20,40,60,80
	Node degree k	5,10,20,40,60
	Mixture coefficient	0,0.25,0.5,0.75,1
Temporal module	Channels	16,32,64
	Dilation	1,2
	kernel numbers	1,2,4,8
	kernel candidates	2,3,4,5,6,7,8,9,10,11,12

which consists of 4 operations which could be spatial module, temporal module, identity module, etc. AutoCTS also has fixed skip connections inside the blocks and among the blocks.

After visualizing above literatures, we propose a disentangled framework as shown in Table 1 and Figure 1(e). We are thus interested in the correlation between the performance and these disjoint factors. Specifically, we consider

- *Spatio-temporal order.* Independent of module specific designs, STGNN models put together spatial/temporal modules in a certain order to process spatial and temporal correlations. In our framework, for a pre-determined number of modules, say 6, each module could either be spatial module (S) or temporal one (T). Thus, it is possible to have T-T-T-T-T-T, i.e. all six temporal modules or vice versa. Note that we do not include parallel connections here, e.g., STSGCN [28] and GMAN [43], because such designs do not perform well, e.g. STFGNN in Table 3.
- *Skip connection.* Earliest works do not consider skip connection [37, 20]. Later works let temporal modules skip to the end and skip internally every two blocks [33]. AutoCTS [32] connects each block to the end and inside the blocks, every later node connects to all previous nodes in a determined way. We consider a much larger skip connection space of over 2 million choices, i.e., each layer could choose to connect to any number of previous layers.
- *Spatial module.* For graph convolution, we consider several different GCNs that are common in Spatio-temporal GNN literature, namely Kipf GCN, Chebyshev GCN and Mixhop GCN. On graph structure, we adopt the setting of AGCRN but with additional kNN for sparsity. To still make use of the provided graph structure, we include a mixture coefficient.
- *Temporal module.* We only consider TCN here since it is easier to be trained than RNN. Specifically, we consider a typical dilated TCN on the time axis of features to extract temporal features. A set of convolution kernels are used and channels of different kernels are aggregated.

We can see that exemplar hand-designed architectures, e.g., STGCN, MTGNN and searched architectures, e.g., AutoCTS, AutoSTG are all included in the above framework. Moreover, it extends the scope of existing works’ to many more possibilities, especially by introducing the space of skip connections and spatio-temporal order.

## 4 Understanding STGNN Hyperparameters and Architectures

In this section, we show how important principles can be observed based on the framework, which help speedup the evaluation and subsequently design better STGNNs. We explain for each factor its motivation, the used methodology and the understanding obtained to answer sequentially the research questions. All experiments in this section are run on two datasets PeMS04 and PeMS08 [13, 28], which are introduced with more details in the Appendix B.

#### 4.1 Understanding training hyperparameters

Different STGNN models naturally need different training hyperparameters. To avoid biased evaluation of models by inheriting hyperparameter setting as existing works, we instead consider a hyperparameter space where bad hyperparameters choices can be removed such that once we run a model multiple times under this space, we obtain a distributionally fair evaluation.

We consider training hyperparameters (HP) that are common from STGNN literature as in Table 2 where the curriculum learning [33, 18] is specified to STGNN here. The current training hyperparameters constitute a space of over  $10^5$  choices and is not practical to efficiently evaluate each model configuration. We aim to independently study the impact by the ranking distribution to remove the hyperparameters that are commonly bad for most architectures. The **ranking distribution** evaluates a model with different hyperparameters to obtain a relative rank and repeats to have distributional patterns. The details of ranking strategy, datasets and choices of hyperparameters are in Appendix B and D.

Table 2: Common training hyperparameters for STGNNs.

Name	Original Scope	Reduced Scope
Learning rate	1e-5, 1e-4, 1e-3, 1e-2, 1e-1	1e-4, 1e-3
Batch size	8,16,32,64,128	8, 32, 128
Optimizer	SGD, Rmsprop, Adam, AdamW, Adamax	Adam, AdamW
Weight decay	0, 1e-1, 1e-2, 1e-3, 1e-4, 1e-5	0, 1e-5
Gradient clip	0, 1, 3, 5, 7, 9	1, 5
Dropout	0, 0.1, 0.3, 0.5, 0.7, 0.9	0, 0.3
Curriculum learning	None, 3, 5, 7	None, 3

**The ranking strategy.** To compare hyperparameter  $P$  of  $K$  choices ( $p_1, p_2, \dots, p_K$ ), we randomly sample a **configuration** under our framework. A configuration contains all necessary choices of hyperparameters and architectures to identify a model. We then replace iteratively the hyperparameter  $P$  of this configuration by all  $K$  choices and train the models one by one. Thus, for each batch of  $K$  runs, we have  $K$  model results whose configurations differ from each other only in the hyperparameter  $P$ . We rank these results ascendingly (smaller error metric means better performance, thus ranks better). We run multiple batches of  $K$  runs and obtain the ranking distribution on each hyperparameter. For example in Figure 2(a), we compare and rank the relative performance of five learning rate choices for a certain model configuration. After many configurations evaluated, we obtain that 1e-3 is the best choice since it ranks first (most frequent) in all runs.

The ranking plot of each hyperparameter is partially given in Figure 2 and the full plots are in Appendix D. The hyperparameters are grouped into three cases.

1. *reduced options*, e.g., learning rate and optimizer could be reduced to a few options;
2. *monotonically related*, e.g., weight decay and dropout rate are (almost) showing a monotonically better as weight decay decreases;
3. *no obvious pattern*, e.g., gradient clip does not show much difference as long as it is activated and same for batch size and curriculum learning.

The training hyperparameter space has been reduced by 500 times as in Table 2, with which we are able to largely increase the efficiency of model evaluation in a fair way.

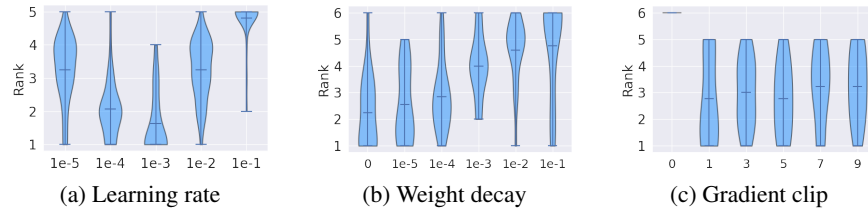


Figure 2: Rankings of three training hyperparameters. A lower rank means better performance. The more a choice ranks first, the better it is in terms of distribution.

## 4.2 Understanding architectures

With such a disentangled framework presented, we aim at the following research questions (RQ) to understand the influence of each factor:

- **RQ1:** What patterns in terms of spatio-temporal order do good models share? How to characterize spatio-temporal order patterns?
- **RQ2:** What patterns in terms of skip connection do good models share? How to characterize skip connection patterns?
- **RQ3:** How do different spatial/temporal module designs influence the performance?
- **RQ4:** How transferable are above patterns in the architecture space?

However, these questions are not trivial to answer. We need to dedicatedly design systematic views over the huge architecture space such that correlations among different architectures and with the final prediction performance can be simultaneously captured. In this way, principles that can indicate goodness of architectures can be generated. Ideally, we want to find cherry regions of architecture space since good architectures should share important designs in common. An architecture in the cherry region has a high probability to achieve good performance. A key challenge in finding such regions lies in the choice of quantitative measures.

### 4.2.1 Spatio-temporal order

The spatio-temporal order specifies the arrangement of both spatial/temporal modules to learn the representation. Intuitively, the representation is influenced by how many spatial/temporal modules we use respectively and how these modules are interleaved to exchange the information. We propose two measures explained below to characterize spatio-temporal order of an architecture:

- **Number of spatial or temporal modules.** Without loss of generality, we observe temporal modules, noted as #T
- **Number of inversed order of a module sequence.** Let the number of inversed order of a pair of modules (either S or T) be 1 if S comes before T, i.e. ST, and 0 otherwise, i.e. TS, SS, TT. The number of inversed order of a module sequence is defined as the sum of that of each pair modules.

The rationale of the above two measures is quantitatively justified in below Remark 1. Specifically, they help us identify a compact triangle region as in Figure 3(a) where different combinatorial choices of S/T orders lie at the corners.

**Remark 1.** We claim that the number of inversed order along with the number of T modules (#T), can well capture the interleaving of a sequence of S-T modules, including the following 4 corner cases: ( $N$  is number of total modules)

- *Case 1: All S modules, #T is 0 and number of inversed order is 0, indicating no temporal feature is needed;*
- *Case 2: All T modules, #T is  $N$  and number of inversed order is 0, indicating no spatial feature is needed;*
- *Case 3:  $\frac{N}{2}$  S modules followed by  $\frac{N}{2}$  T modules, #T is  $\frac{N}{2}$  and number of inversed order is 0, indicating the spatial information should strictly be processed before temporal information;*
- *Case 4:  $\frac{N}{2}$  T modules followed by  $\frac{N}{2}$  S modules, #T is  $\frac{N}{2}$  and number of inversed order is  $\frac{N^2}{4}$ , indicating the temporal information should strictly be processed before spatial information.*

We fix the framework factors except for the spatio-temporal order. As mentioned above, we relax the spatio-temporal order under 6 temporal and/or spatial modules, leading to a space of  $2^6 = 64$  choices. We use the above mentioned measures to observe the order. Results are in Figure 3. We can see that for spatio-temporal order, all corner cases perform terribly. A clear region for 2 to 5 T blocks and for number of inversions between 3 and 8 exists and shows superior performance. This suggests that (1) we might have good performance for many temporal modules but too many spatial modules is not a good idea. This could be because of the oversmoothing phenomenon observed in GNN [35] but further study is required; (2) a large number of inversion usually returns better performance which means that both modules are encouraged to interleave more frequently.

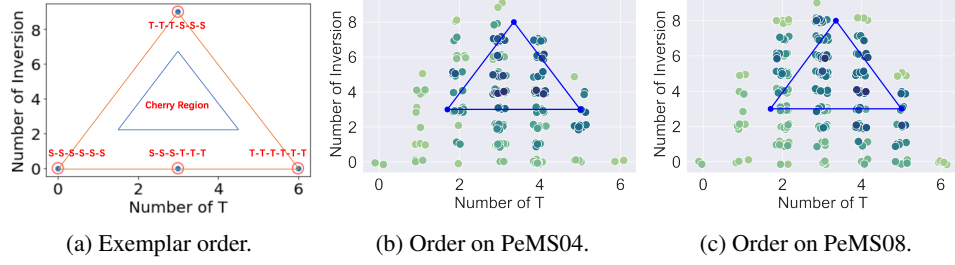


Figure 3: Spatio-temporal order understanding on different datasets. For (b)(c), we normalize performances to show them in a colormap. The darker the circle, the better the performance.

#### 4.2.2 Skip connection

As the representation flows from the input to the output without a loop, we are motivated to formulate skip connections in STGNNs as a Directed Acyclic Graph (DAG) with a straight though flow path connecting sequentially each module (Figure 1(e)). Afterwards, we need to consider proper measures to observe the skip space as it is large. Intuitively, the performance of a certain skip pattern depends on how many we skip and where we skip. The former can be measured by number of paths from input module to output module and the latter can be measured by average shortest length for each pair of modules, both defined below. We propose another two measures explained below to characterize skip connection of an architecture:

- **Number of path**, noted as (#Path). It is defined as the total number of trajectories from the predetermined input node to the predetermined output node.
- **Average Shortest Length**, noted as ASL. It is defined as the average of shortest lengths from non-output nodes to the output node predetermined.

The rationale of the above two measures is quantitatively justified in below Remark 2. Specifically, they help us identify a compact right triangle region as in Figure 4(a) where different combinatorial choices of skip connections lie at the corners.

**Remark 2.** We claim that Average Shortest Length (ASL) and number of path(#Path) in a Directed Acyclic Graph can well capture the spatio-temporal information flow, including the following 3 corner cases:

- *Case 1: ASL is low and #Path is low, i.e., limited number of skips which mostly connects to the output node, where most literature falls in;*
- *Case 2: ASL is low and #Path is high, i.e., a lot of skip connections among which many skips to the output node, potentially obfuscating the propagated message;*
- *Case 3: ASL is high and #Path is low, i.e., almost no additional skips, which is the case of earliest STGCN model [37].*

In Figure 4, we show the skip experiment on two datasets and on two spatio-temporal orders that have been used in the literature. Details are given in Appendix D. All these four experiments follow a similar pattern. First, all corner cases perform badly. In the center rectangle, we observe a clear region where the performance are much better, especially for ASL between 1.1 to 1.8 and #Path between 10 and 40. Notably, this region exists for different datasets and for different spatio-temporal orders, showing its universality in terms of principles. This observation can be exemplified by other works e.g. residual connection [14], JK-Net [35], message passing [2]. Especially for GNN, we want to reinforce the remote message, e.g. the input node, in later nodes, thus ASL will be low. We also do not want to have too many messages that obfuscate the valuable message, thus #Path cannot be too high.

#### 4.2.3 Spatial/temporal module design

Different choices of spatial/temporal module influence the way the spatial/temporal representation are processed. In the literature, diverse modifications on TCN and GCN have been proposed. We use the same ranking distribution as in understanding training hyperparameters to study all the module

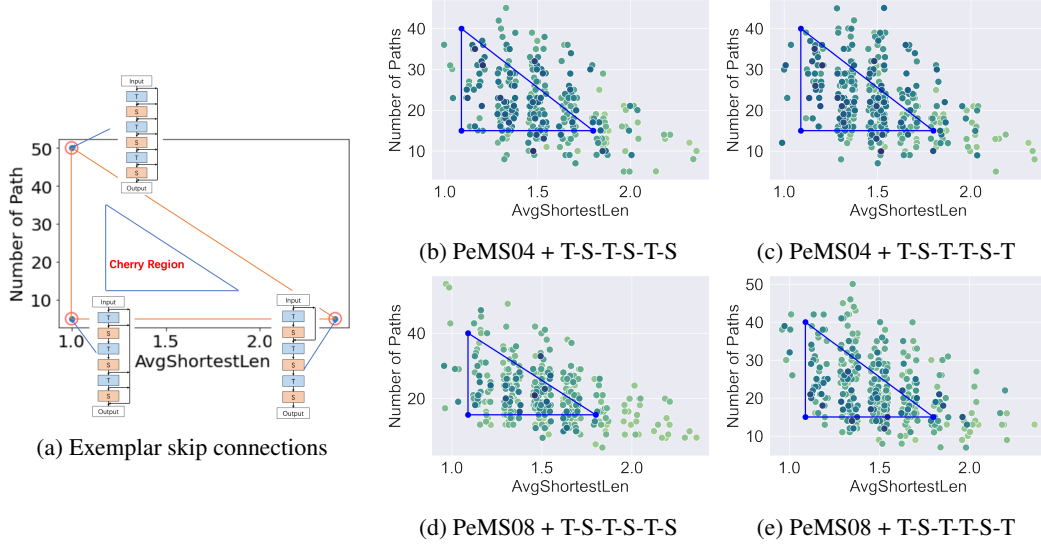


Figure 4: Skip understanding on different spatio-temporal orders and different datasets. For (b)(c)(d)(e), we normalize performances to show them in a colormap. The darker the circle, the better the performance.

270 choices. The ranking plots are partially shown in Figure 5 and full plot is given in Appendix D. We  
 271 remove a few bad options according to the ranking, e.g., Cheb GCN, node degree 5, etc. Then, we  
 272 also find we could largely reduce number of parameters by reducing channels temporarily due to a  
 273 high performance correlation, as show in Appendix D.

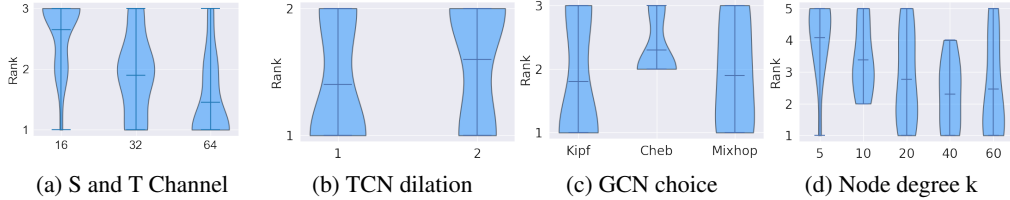


Figure 5: Rankings of several design choices. A lower rank means better performance.

#### 274 4.2.4 Summary of designing principles

275 With the empirical results in this section, after refining the training hyperparameter space, we find that  
 276 good architectures share similar patterns. Firstly, Two cherry regions exist as shown in Figure 3(a)  
 277 and Figure 4(a), answering **RQ1** and **RQ2**. Concretely, we tend to consider architectures which  
 278 have appropriate number of temporal modules, inversion number, averaged shortest length, and  
 279 finally number of paths. The appropriate scope can be interpreted from the view of spatio-temporal  
 280 information flow and message passing with jumping information. For module designs, we remove  
 281 some bad choices as by ranking and discover that we could largely reduce number of parameters by  
 282 reducing channel size due to high performance correlation, answering **RQ3**. All the experiments in  
 283 this section are conducted on two datasets PeMS04 and PeMS08, showing the generability of the  
 284 discoveries, thus answering **RQ4**.

## 285 5 Searching Better STGNN Models in a Simplified Way

### 286 5.1 A simplified but strong NAS baseline

287 Here, to give an example on how to leverage these understandings, we show a simple but strong  
 288 NAS baseline based on Random Search (RS) that is able to find better STGNN models efficiently,  
 289 titled **SimpleSTG** as illustrated in Figure 6. We use random search with the cherry regions found



by understanding results on typical datasets. Specifically, we evaluate a *random* sample if it lies in regions, otherwise reject it and re-sample. Note that more complex search strategies leveraging our understanding surely exists and probably will be more efficient than our simple baseline, e.g. Evolutionary, Bayesian or even differentiable. But our goal is to demonstrate how to utilize our understanding and how effective it could be to design architectures with better understandings. We stick to this baseline for further comparison.

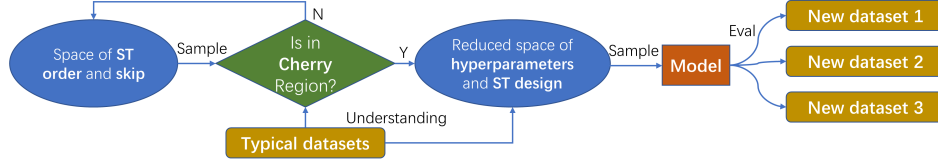


Figure 6: Illustration of the pipeline in our simple random search baseline SimpleSTG.

## 5.2 Overall performance comparison

In this part, we compare more thoroughly our models with related works. From Table 3, incorporating architecture priors, the proposed simple baseline could find novel STGNN models better than hand-designed and NAS-based methods. The other baselines are introduced as follows. Vector Auto-Regression (VAR) is a baseline often for sanity check.

ASTGCN [13] adopts attention mechanism for both spatial and temporal modeling. STSGCN [28] proposes synchronous modules to leverage the heterogeneities in spatial-temporal data. STSGCN [28] stacks multiple localized GCN for extraction of correlations. STGODE [10] uses a tensor-based ordinary differential equation. Z-GCNETs [5] is GRU-based and introduces a time zigzag persistence to integrate with GCN. Others have been discussed in Section 1. Note that there are many NAS methods for pure GNN such as [21, 42]. These methods cannot be used directly in STGNN because of the difference in search space.

Table 3: Prediction performance comparison of different models. Bold number denotes the best and underscored number denotes the second best. Numbers in the parentheses indicate the standard deviation.

Model	PeMS03		PeMS04		PeMS07		PeMS08	
	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE
VAR	23.65(0.00)	38.26(0.00)	23.75(0.00)	36.66(0.00)	75.63(0.00)	115.2(0.00)	23.46(0.00)	36.33(0.00)
STGCN	17.49(0.46)	30.12(0.70)	22.70(0.64)	35.55(0.75)	25.38(0.49)	38.78(0.58)	18.02(0.14)	27.83(0.20)
DCRNN	18.18(0.15)	30.31(0.25)	24.70(0.22)	38.12(0.26)	25.30(0.52)	38.58(0.70)	17.86(0.03)	27.83(0.05)
ASTGCN	17.69(1.43)	29.66(1.68)	22.93(1.29)	35.22(1.90)	28.05(2.34)	42.57(3.31)	18.61(0.40)	28.16(0.48)
STSGCN	17.48(0.15)	29.21(0.56)	21.19(0.10)	33.65(0.20)	24.26(0.14)	39.03(0.27)	17.13(0.09)	26.80(0.18)
AGCRN	16.58(0.10)	27.48(0.14)	19.83(0.06)	32.26(0.12)	24.21(0.21)	37.66(0.24)	15.95(0.18)	25.22(0.22)
MTGNN	15.23(0.03)	26.12(0.20)	19.25(0.03)	31.65(0.21)	21.28(0.11)	34.31(0.16)	15.86(0.10)	24.93(0.19)
STGODE	-	-	20.84(0.00)	32.82(0.00)	-	-	16.81(0.00)	25.97(0.00)
Z-GCNETs	-	-	19.90(0.00)	32.66(0.00)	-	-	16.12(0.00)	25.74(0.00)
STFGNN	16.77(0.09)	28.34(0.46)	19.83(0.06)	31.88(0.14)	22.07(0.11)	35.80(0.18)	16.64(0.09)	26.22(0.15)
AutoCTS	<u>14.71</u> (0.40)	<u>24.54</u> (0.33)	<u>19.13</u> (0.21)	<b>30.44</b> (0.24)	<u>20.93</u> (0.35)	<u>33.69</u> (0.29)	<u>14.82</u> (0.17)	<b>23.64</b> (0.10)
<b>SimpleSTG (ours)</b>	<b>14.45</b> (0.10)	<b>24.35</b> (0.13)	<b>18.56</b> (0.26)	<u>30.71</u> (0.41)	<b>19.80</b> (0.20)	<b>33.03</b> (0.18)	<b>14.64</b> (0.10)	<u>23.77</u> (0.12)

The **generalizability** of the distilled principles is demonstrated in two ways. First, in Section 4, we show at the same time empirical results on two datasets and similar principles are observed. Second, in Table 3, note that we search only one model and evaluate this model on more datasets covering different regions of California (Appx B) instead of searching the best model per dataset as in other NAS methods. We show further the comparison on a new and different dataset NE-BJ released by [18]. The NE-BJ dataset contains traffic information in Beijing, which is totally different from commonly used California datasets. The metric is MAE and RMSE on 15 mins, 30 mins and 60 mins prediction to comprehensively evaluate the effectiveness of our model. In almost all test cases, our

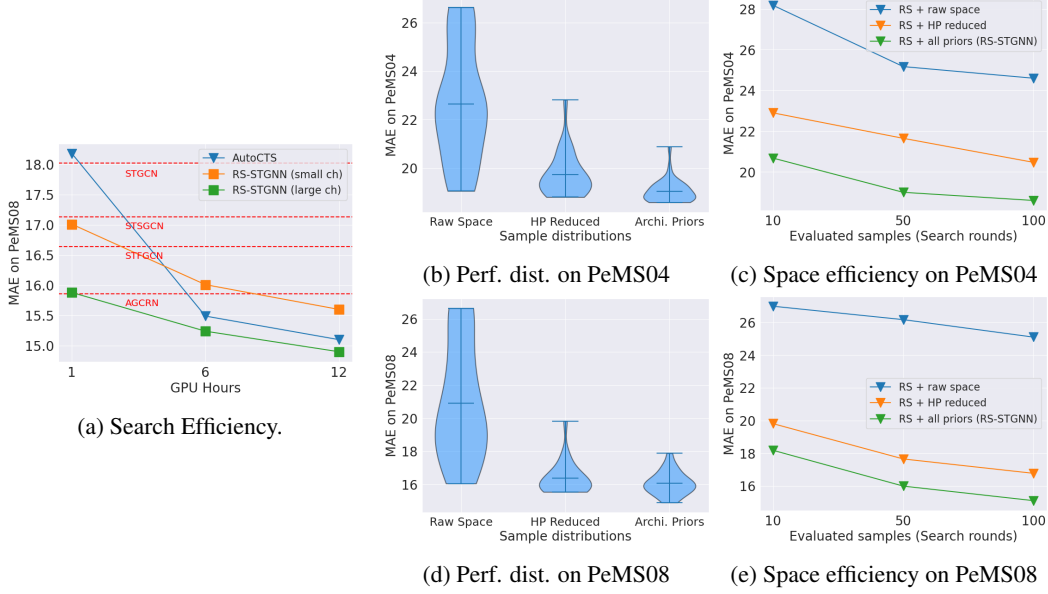


Figure 7: Empirical results on the search efficiency and effectiveness of space.

baseline outperforms by a large margin. We also provide evaluation on another non-traffic dataset and the configuration of searched model in Appx E.

Table 4: Performance comparison on NE-BJ. Bold number denotes the best and underscored number denotes the second best. Instead of averaged 12 steps, we show results on 3/6/12 steps each for thoroughness.

Model	15 mins		30 mins		60 mins	
	MAE	RMSE	MAE	RMSE	MAE	RMSE
VAR	5.42	8.16	5.76	9.07	6.14	9.65
DCRNN	3.84	6.84	4.51	8.49	5.15	9.77
STGCN	5.02	8.34	5.10	8.55	5.39	<b>9.09</b>
ASTGCN	4.43	7.34	5.31	8.86	6.29	10.31
AGCRN	3.90	6.81	4.55	<u>8.32</u>	5.06	9.54
GMAN	4.08	7.63	<u>4.42</u>	<u>8.45</u>	<b>4.80</b>	9.18
AutoCTS	3.91	<u>6.70</u>	4.69	<b>8.21</b>	5.64	9.80
<b>SimpleSTG (ours)</b>	<b>3.71</b>	<b>6.69</b>	<b>4.33</b>	<b>8.21</b>	<u>4.85</u>	9.34

### 5.3 Search efficiency and effectiveness of space

We further demonstrate the effectiveness of our understanding on training hyperparameters and architectures in Figure 7(a-e) on two datasets. In the Figure 7(a), we evaluate the search efficiency in terms of GPU hours and compare NAS methods with hand designed results on dataset PeMS08. The orange line is our accelerated search in small channel and the green line is the same architecture as in orange line but with large channel size. The motivation and correlation of performances with different channel sizes are illustrated in Appx D. It can be found that our NAS method, though implemented with a simple random search strategy, can be very efficient in terms of GPU hours. Note also that our random search baseline is straightforward to adapt on distributed systems without communication head while AutoCTS cannot easily achieve significant acceleration. In Figure 7(b)(d), we show the samples' performances in different spaces. With architecture priors and reduced hyperparameter space, the performance distribution is significantly better. In Figure 7(c)(e), we show the random search efficiency considering different priors. Both plots show the effectiveness of our understanding and concluded principles.

## 6 Conclusion and discussion

In this work, we revisit Spatio-Temporal Graph Neural Networks in traffic prediction to study how to properly design models in a principled way. We propose a framework that disentangles the choices into three groups: spatial/temporal module designs, spatio-temporal order and skip connection. We understand qualitatively and quantitatively the influence of each disjoint factor and conclude the principles behind the architecture design. To illustrate how these understandings could help find better STGNN models, we propose a simple strategy to efficiently and effectively search models that outperform all other works.

## References

- [1] Lei Bai, Lina Yao, Can Li, Xianzhi Wang, and Can Wang. Adaptive graph convolutional recurrent network for traffic forecasting. In *Advances in Neural Information Processing Systems*, 2020.
- [2] Peter W. Battaglia, Jessica B. Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinícius Flores Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, Çağlar Gülçehre, H. Francis Song, Andrew J. Ballard, Justin Gilmer, George E. Dahl, Ashish Vaswani, Kelsey R. Allen, Charles Nash, Victoria Langston, Chris Dyer, Nicolas Heess, Daan Wierstra, Pushmeet Kohli, Matthew M. Botvinick, Oriol Vinyals, Yujia Li, and Razvan Pascanu. Relational inductive biases, deep learning, and graph networks. *CoRR*, 2018.
- [3] Irwan Bello, William Fedus, Xianzhi Du, Ekin Dogus Cubuk, Aravind Srinivas, Tsung-Yi Lin, Jonathon Shlens, and Barret Zoph. Revisiting resnets: Improved training and scaling strategies. In *Advances in Neural Information Processing Systems*, 2021.
- [4] Xiangning Chen, Ruochen Wang, Minhao Cheng, Xiaocheng Tang, and Cho-Jui Hsieh. Drnas: Dirichlet neural architecture search. In *International Conference on Learning Representations*, 2021.
- [5] Yuzhou Chen, Ignacio Segovia-Dominguez, and Yulia R. Gel. Z-gcnets: Time zigzags at graph convolutional networks for time series forecasting. In *International Conference on Machine Learning*, 2021.
- [6] Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural machine translation: Encoder-decoder approaches. In *Eighth Workshop on Syntax, Semantics and Structure in Statistical Translation, SSST@EMNLP*. Association for Computational Linguistics, 2014.
- [7] Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Twins: Revisiting the design of spatial attention in vision transformers. In *Advances in Neural Information Processing Systems*, 2021.
- [8] Zhiyong Cui, Kristian Henrickson, Ruimin Ke, and Yin Hai Wang. Traffic graph convolutional recurrent neural network: A deep learning framework for network-scale traffic learning and forecasting. *IEEE Transactions on Intelligent Transportation Systems*, 2020.
- [9] Moshe Eliasof, Eldad Haber, and Eran Treister. PDE-GCN: novel architectures for graph neural networks motivated by partial differential equations. In *Advances in Neural Information Processing Systems*, pages 3836–3849, 2021.
- [10] Zheng Fang, Qingqing Long, Guojie Song, and Kunqing Xie. Spatial-temporal graph ODE networks for traffic flow forecasting. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2021.
- [11] Rui Fu, Zuo Zhang, and Li Li. Using lstm and gru neural network methods for traffic flow prediction. In *Youth Academic Annual Conference of Chinese Association of Automation*, 2016.
- [12] Yury Gorishniy, Ivan Rubachev, Valentin Khrulkov, and Artem Babenko. Revisiting deep learning models for tabular data. In *Advances in Neural Information Processing Systems*, 2021.

- [13] Shengnan Guo, Youfang Lin, Ning Feng, Chao Song, and Huaiyu Wan. Attention based spatial-temporal graph convolutional networks for traffic flow forecasting. In *Association for the Advancement of Artificial Intelligence*, 2019.
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Computer Vision and Pattern Recognition*, 2016.
- [15] Renhe Jiang, Du Yin, Zhaonan Wang, Yizhuo Wang, Jiewen Deng, Hangchen Liu, Zekun Cai, Jinliang Deng, Xuan Song, and Ryosuke Shibasaki. DI-traff: Survey and benchmark of deep learning models for urban traffic prediction. In *ACM International Conference on Information and Knowledge Management*, 2021.
- [16] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017.
- [17] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson, Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropagation applied to handwritten zip code recognition. *Neural computation*, 1989.
- [18] Fuxian Li, Jie Feng, Huan Yan, Guangyin Jin, Depeng Jin, and Yong Li. Dynamic graph convolutional recurrent network for traffic prediction: Benchmark and solution. *CoRR*, 2021.
- [19] Mengzhang Li and Zhanxing Zhu. Spatial-temporal fusion graph neural networks for traffic flow forecasting. In *AAAI Conference on Artificial Intelligence*, 2021.
- [20] Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. In *International Conference on Learning Representations*, 2018.
- [21] Yanxi Li, Zean Wen, Yunhe Wang, and Chang Xu. One-shot graph neural architecture search with dynamic search space. In *AAAI Conference on Artificial Intelligence*, pages 8510–8517, 2021.
- [22] Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: differentiable architecture search. In *International Conference on Learning Representations*, 2019.
- [23] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.
- [24] Attila M. Nagy and Vilmos Simon. Survey on traffic prediction in smart cities. *Pervasive and Mobile Computing*, 2018.
- [25] Zheyi Pan, Songyu Ke, Xiaodu Yang, Yuxuan Liang, Yong Yu, Junbo Zhang, and Yu Zheng. AutoSTG: Neural architecture search for predictions of spatio-temporal graph. In *The Web Conference*, 2021.
- [26] Bin Ran and David Boyce. *Modeling dynamic transportation networks: an intelligent transportation system oriented approach*. Springer Science & Business Media, 2012.
- [27] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks and Learning Systems*, 2009.
- [28] Chao Song, Youfang Lin, Shengnan Guo, and Huaiyu Wan. Spatial-temporal synchronous graph convolutional networks: A new framework for spatial-temporal network data forecasting. In *Association for the Advancement of Artificial Intelligence*, 2020.
- [29] Mingxing Tan and Quoc V. Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, 2019.
- [30] Leye Wang, Di Chai, Xuanzhe Liu, Liyue Chen, and Kai Chen. Exploring the generalizability of spatio-temporal traffic prediction: Meta-modeling and an analytic framework. *IEEE Transactions on Knowledge and Data Engineering*, 2021.

- 425 [31] Colin White, Arber Zela, Robin Ru, Yang Liu, and Frank Hutter. How powerful are performance  
426 predictors in neural architecture search? In *Advances in Neural Information Processing Systems*,  
427 pages 28454–28469, 2021.
- 428 [32] Xinle Wu, Dalin Zhang, Chenjuan Guo, Chaoyang He, Bin Yang, and Christian S. Jensen.  
429 Autocts: Automated correlated time series forecasting. *International Conference on Very Large*  
430 *Data Bases*, 2021.
- 431 [33] Zonghan Wu, Shirui Pan, Guodong Long, Jing Jiang, Xiaojun Chang, and Chengqi Zhang.  
432 Connecting the dots: Multivariate time series forecasting with graph neural networks. In *ACM*  
433 *SIGKDD Conference on Knowledge Discovery and Data Mining*, 2020.
- 434 [34] Sirui Xie, Hehui Zheng, Chunxiao Liu, and Liang Lin. SNAS: stochastic neural architecture  
435 search. In *International Conference on Learning Representations*. OpenReview.net, 2019.
- 436 [35] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and  
437 Stefanie Jegelka. Representation learning on graphs with jumping knowledge networks. In  
438 *International Conference on Machine Learning*, 2018.
- 439 [36] Jiaxuan You, Zhitao Ying, and Jure Leskovec. Design space for graph neural networks. In  
440 *Advances in Neural Information Processing Systems*, 2020.
- 441 [37] Bing Yu, Haoteng Yin, and Zhanxing Zhu. Spatio-temporal graph convolutional networks: A  
442 deep learning framework for traffic forecasting. In *International Joint Conference on Artificial*  
443 *Intelligence*, 2018.
- 444 [38] Haiyang Yu, Zhihai Wu, Shuqin Wang, Yunpeng Wang, and Xiaolei Ma. Spatiotemporal  
445 recurrent convolutional networks for traffic prediction in transportation networks. *Sensors*,  
446 2017.
- 447 [39] H. Zare Moayed and M.A. Masnadi-Shirazi. Arima model for network traffic prediction and  
448 anomaly detection. In *International Symposium on Information Technology*, 2008.
- 449 [40] Junping Zhang, Fei-Yue Wang, Kunfeng Wang, Wei-Hua Lin, Xin Xu, and Cheng Chen.  
450 Data-driven intelligent transportation systems: A survey. *IEEE Transactions on Intelligent*  
451 *Transportation Systems*, 2011.
- 452 [41] Yongqi Zhang, Zhanke Zhou, Quanming Yao, and Yong Li. Efficient hyper-parameter search  
453 for knowledge graph embedding. In *Annual Meeting of the Association for Computational*  
454 *Linguistics*, pages 2715–2735, 2022.
- 455 [42] Huan Zhao, Quanming Yao, and Weiwei Tu. Search to aggregate neighborhood for graph neural  
456 network. In *IEEE International Conference on Data Engineering*, pages 552–563, 2021.
- 457 [43] Chuanpan Zheng, Xiaoliang Fan, Cheng Wang, and Jianzhong Qi. GMAN: A graph multi-  
458 attention network for traffic prediction. In *AAAI Conference on Artificial Intelligence*, 2020.

## A Related works

Our proposed framework differs from existing works as follows. DL-Traff [15] surveys the literature and proposes a framework based on the modeling method on spatial axis and temporal axis (e.g., CNN, GRU, GNN, etc.). DGCRN [18] considers additional spatial topology construction and external features to classify different methods. STMeta [30] proposes another analytical framework concerning spatio-temporal information modeling. Our framework differs from then in that we also consider macroscopic factors and training hyperparameters. The purpose of our framework is also different because we use our framework to understand the architecture space which potentially facilitates further architecture search. On the other hand, AutoCTS [32] and AutoSTG [25] proposes their own framework of STGNN to search architectures, but they do not consider skip connection or training hyperparameters. AutoCTS and AutoSTG do not try to understand the space.

Our methodologies of understanding have been inspired by many other works, including architecture understanding in Convolution Neural Networks, Graph Neural Networks and Neural Architecture Search [36], different communities’ common discovery that training hyperparameters once well tuned can boost the model performance by a large margin and probably influence the models’ evaluation [29, 36]. Previous works on surveys, benchmarks, and NAS of traffic prediction should also be credited, as discussed in the previous paragraph.

## B Details on the datasets and problem setting

### B.1 Dataset

To evaluate quantitatively the performance of different models, we experiment on four public real world datasets: PeMS03, PeMS04, PeMS07 and PeMS08 [13, 28, 1, 10]. These datasets can be accessed on GitHub<sup>2</sup>. These datasets come from Caltrans Performance Measurement System (PeMS) and contain the measures of highway traffic flow in different regions. The raw temporal signal was recorded every 30 seconds and we aggregate both datasets into 5-mins regular traffic measures the same way as previous works. Detailed statistics are given in Table 5. Datasets are splitted in a 6,2,2 manner. A standard z-score normalization is applied which subtracts and divides data by mean and standard deviation of training split. We DO NOT include additional hand-craft features e.g. calendar, holiday or time of the day. In all our experiments of Section 4, we use both PeMS04 and PeMS08 for understanding general principles. In Section 5, we show overall model performance comparisons on all these four datasets.

Table 5: Dataset statistics

Dataset	Region	#Timestamps	#Sensors	Time Period
PeMS03	Marysville	26208	358	01/09/2018 - 30/11/2018
PeMS04	Oakland	16992	307	01/01/2018 - 28/02/2018
PeMS07	Los Angeles	28224	883	01/05/2017 - 31/08/2017
PeMS08	San Bernardino	17856	170	01/07/2016 - 31/08/2016

### B.2 Problem setting

In Section 2, we give a formal problem definition on spatio-temporal traffic prediction. In practice, following the literatue [38, 8, 1, 33], we set  $P = Q = 12$ , i.e., predicting the next one hour based on past hour. Concretely, if we predict on dataset PeMS04 from Table 5, the input tensor is of shape: (307, 12, 1) and the output tensor is of shape: (307, 12, 1). 1 is the channels of traffic signal or feature. Since we do not include additional calendar or hand-designed features, this channel remains 1.

<sup>2</sup><https://github.com/Davidham3/ASTGCN/tree/master/data>

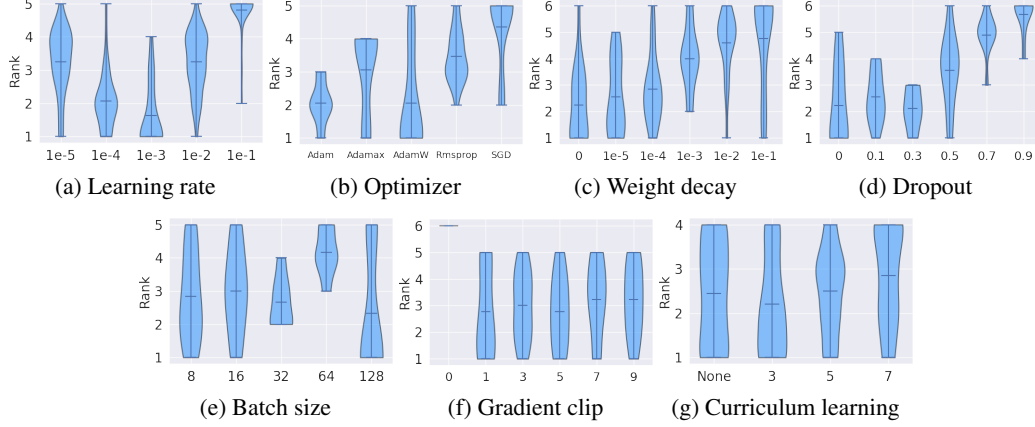


Figure 8: Rankings of different training hyperparameters. A lower rank means better performance. The more a choice ranks first, the better it is in terms of distribution.

## C Details on the disentangled framework

### C.1 Curriculum learning

As pointed out by [33, 18], curriculum learning improves consistently the multi-step prediction. To use curriculum learning in traffic prediction, we first start with one step prediction with all training data (easiest task) and gradually increase the prediction horizon to deal with more difficult task. This curriculum learning strategy brings another hyperparameter: number of epochs to increase the task difficulty.

### C.2 Calculation of skip connection choices

We have in total  $2^{\frac{(L-1)(L-2)}{2}}$  skip connection choices where  $L$  is number of modules (including input module, spatio-temporal modules, and output module). The first module has  $(L - 2)$  skip placeholders. The second module has  $(L - 3)$  skip placeholders, and so on. In the case of 8 modules (6 spatio-temporal module plus 1 input and 1 output module), we will have a space of  $2^{6+5+4+3+2+1} = 2^{21} = 2.09$  millions choices.

## D Details on the understanding results

### D.1 More observations on training hyperparameters

Additional interesting observations could be made in Figure 2. For example, batch size 64 which most literature uses is worst in our ranking. The optimizer AdamW [23] is not considered by any literature in traffic prediction but showed best in our rank. All these give evidences to certain biases in training hyperparameters. As a result, to compare fairly different model architectures, we need to take into account the training hyperparameters.

### D.2 More on spatial/temporal module design

Different choices of spatial/temporal module influence the way the spatial/temporal representation are processed. In the literature, diverse modifications on TCN and GCN have been proposed. With limited ablation study, it is hard to conclude which part really contributes to the performance gain, thus hinder our understanding on spatio-temporal modeling.

We use the same ranking distribution as in understanding training hyperparameters to study all the module choices. Due to the huge impact on model parameters, we study in addition the channel correlation whose detailed plots are in Appendix D.

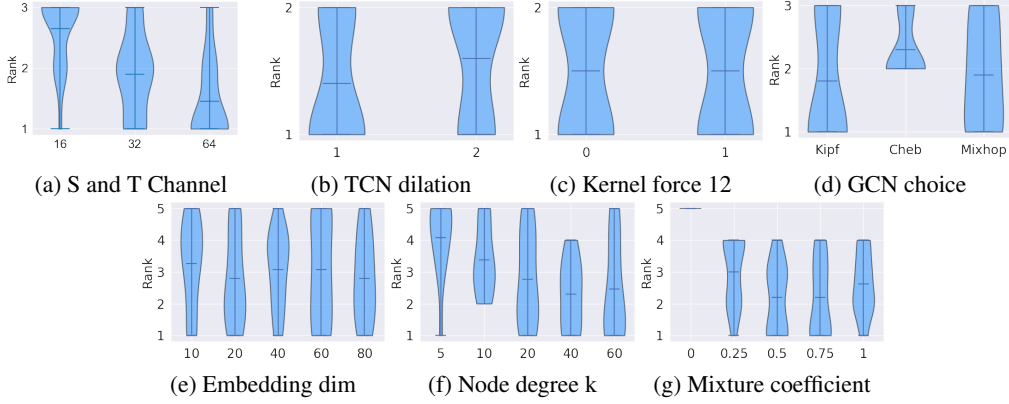


Figure 9: Rankings of different design choices. A lower rank means better performance.

523 The full ranking plots are shown in Figure 9. Firstly, except for spatial/temporal channels, most  
 524 designs from the literature are useful. We remove a few bad options according to the ranking, e.g.,  
 525 Cheb GCN, zero mixture coefficient, etc. Then, spatial/temporal channel influences a lot the trainable  
 526 parameters and shows consistent goodness as in Figure 9(a), we study further the performance  
 527 correlation between channel 64 and channel 32. We find that channel 32 reaches a correlation up to  
 528 0.85 with channel 64 while having only half of the parameters. Thus, we use channel 32 for all the  
 529 intermediate study for acceleration and change back to 64 for final evaluation.

### 530 D.3 More on channel correlation

531 We have found in Section 4 that channel rank is very significant, i.e., larger channel gives better  
 532 performance. Actually, it is usually believed that larger channels can perform better but also bring  
 533 more trainable parameters [29]. Thus, after ranking them, we consider in addition to study the  
 534 performance correlation between different channel choices to understand its impact and further opens  
 535 the possibility to scale down models to accelerate the execution without messing up models' relative  
 536 performance. We choose spatial module channel to be the same as temporal module channel for  
 537 simplicity.

538 First, the channel correlation can be visualized in Figure 10. We show that on both datasets (PeMS04  
 539 and PeMS08), channel 32 achieves a high performance correlation with channel 64 while the  
 540 correlation between channel 16 and channel 64 are not stable. On the other hand, channel 32 already  
 541 reduces half of the trainable parameters compared to channel 64 and significantly lowers the training  
 542 time cost. As a result, we will use channel 32 for the understanding study except for the final  
 543 performance comparison where we use channel 64.

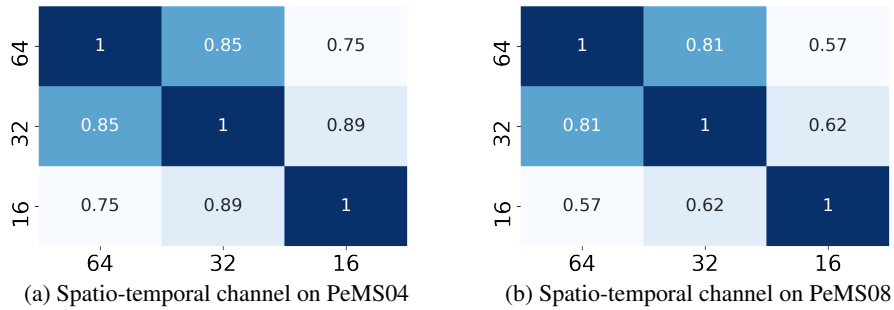


Figure 10: Performance correlation with different channels.



#### 544 D.4 More on the chosen measures

545 Note that other measures exist but may not be the most appropriate in our study. In the example of  
546 skip connection, shortest length from input module to output module reflects also where we skip. But  
547 we can only distinguish the case of shortest length 1 and the other (because we have too many cases  
548 that coincidence on shortest length 2) which is not informative enough.

#### 549 D.5 More on the skip connection experiments

550 To better understand the impact of skip patterns, we fix the other parts of the framework. In this part,  
551 we study the case of 6 spatial/temporal modules thus  $N$  is 8. As by calculation in Appendix C, we  
552 explore the skip space of  $2^{21} = 2907152$  choices. For a certain dataset and a certain fixed the other  
553 configuration, we randomly sample among these two million choices and record their performance  
554 with the above mentioned measures.

### 555 E Details on the searched model

#### 556 E.1 Configuration of searched model

557 The found best model on PeMS08 uses the following configuration. Learning rate 1e-3; Batch size  
558 128; Optimizer AdamW; Weight decay 0; Gradient clip 5; Dropout 0.3; Curriculum learning 3 or  
559 5; Temporal channels 64; Dilation 1; Kernel set 7,8; GCN channels 64; Graph convolution Mixhop  
560 GCN; Embedding dim 20; Node degree 40; Mixture coefficient 0.75.

#### 561 E.2 Evaluation metric

562 Following the literature [38, 8, 1, 33], we use past 12 steps to predict next 12 steps. Both metrics  
563 MAE and RMSE are calculated as the average over future 12 steps. Each step is 5min so we are  
564 predicting one hour ahead using past hour. In Sec 5, we also give results on 3/6/12 steps (15/30/60  
565 mins) to compare more thoroughly.

#### 566 E.3 More on the compared baselines

567 **MTGNN.** When running MTGNN, we notice that it applies a step-number-based curriculum learning.  
568 If we apply directly on other datasets of different size, it will not be able to train on all 12 prediction  
569 steps and perform poorly. As a result, we increase the epochs to 200 instead of 100 to have MTGNN  
570 trained on full 12 horizons. We leave other hyperparameters as default. Our executed results are close  
571 to the reproduced ones from AutoCTS [32].

572 **AutoSTG.** We find two works of AutoML applied to STGNN modeling, AutoCTS and AutoSTG.  
573 We compare only with AutoCTR because AutoSTG requires additional meta information from  
574 geographic database to construct their model which largely limits their applicability, as noted by  
575 AutoCTS [32].

#### 576 E.4 Parameter comparison

577 During our study, we also notice that the searched models are larger in terms of trainable parameters.  
578 In the literature, STGNN models are often smaller than 1 million(M) parameters while our searched  
579 model has parameters of 1.7M. To compare more fairly and understand better the models, we  
580 experiment further and show two things (1) more parameters do not always help (2) by compressing  
581 our searched model to same level of parameters with literature, we can still outperform them indicating  
582 our other choices of architecture indeed help. The results is summarized in Table 6.

Table 6: Prediction performance comparison with models of different scale (Dataset PeMS08).

Model	#Param	MAE	RMSE
AGCRN	0.8M	15.95	25.22
AGCRN (large)	1.8M	16.56	26.43
MTGNN	0.4M	15.86	24.72
MTGNN (large)	1.9M	16.02	25.30
Our searched model (ch=16)	0.4M	15.13	24.13
Our searched model (ch=32)	0.8M	15.04	24.03
Our searched model (ch=64)	1.7M	14.75	23.82

## 583 E.5 Evaluation on non-traffic dataset

584 To even further evaluate the generalization of our method, we use another non-traffic dataset “Elec-  
585 tricity” [32] and measure on 3/6/12 steps as well as average over 12 steps.

Table 7: Performance comparison on Electricity.

Model	Average		3 steps		6 steps		12 steps	
	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE
STGCN	423.1	3923.9	350.5	3706.5	414.3	3918.4	479.8	4111.2
DCRNN	400.5	3640.3	320.6	3356.5	380.8	3634.3	461.2	3919.6
AGCRN	297.1	2353.1	274.6	2066.6	304.3	2329.4	312.4	2475.4
MTGNN	261.0	2624.4	238.6	2331.3	267.9	2631.5	283.9	2867.7
AutoCTS	255.4	2101.6	230.9	1986.5	258.4	2115.9	277.1	2256.3
<b>SimpleSTG (ours)</b>	<b>222.9</b>	<b>1926.9</b>	<b>214.8</b>	<b>1835.7</b>	<b>218.5</b>	<b>1843.7</b>	<b>238.0</b>	<b>2083.6</b>
Improvement(%)	12.7	8.3	7	7.6	15.4	12.9	14.1	7.7

## NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and precede the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

**The checklist answers are an integral part of your paper submission.** They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

IMPORTANT, please:

- **Delete this instruction block, but keep the section heading “NeurIPS paper checklist”,**
- **Keep the checklist subsection headings, questions/answers and guidelines below.**
- **Do not modify the questions and only use the provided macros for your answers.**

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope?

Answer: [Yes]

Justification: Our main claims are listed in section 1 and highlight this paper’s contributions.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We have discussed the potential limitations in section 6, including the application domain and considered temporal modeling methods.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: We have discussed the assumptions in the section 4 when we introduce graph characteristics.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We discuss the experiment settings in section 4, section 5 and appendix B, C, D, E including dataset, baseline, implementation and results. We also include our code in the github link <https://github.com/AutoML-Research/SimpleSTG>.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 5. Open Access to Data and Code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: All the data and code are open access. The dataset information is mentioned in section 4, 5, Appendix B and the code is available in github <https://github.com/AutoML-Research/SimpleSTG>.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.

- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [\[Yes\]](#)

Justification: Hyperparameter sampling is one of our method designs as explained in section 4 and 5. We also specify the training and test details in the same sections and appendix B.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [\[Yes\]](#)

Justification: Our main experiments in section 4 and 5 all include variance as error bars. The experiments are based on multiple runs with mean and variance calculated.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [\[Yes\]](#)

Justification: In section 5, we reported the GPU hours that are required to generate reasonable answers w.r.t. multiple baselines.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: We review the code of ethics and our paper conform with it.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We discuss the broader impacts in the section "Broader Impact Statement".

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [Yes]

Justification: We have discussed the details of the used datasets and models in section 4, 5 and appendix B.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for Existing Assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We explicitly mentioned and cited the source of the data used in our paper [13, 28, 1, 10]. These datasets can be accessed on GitHub <https://github.com/Davidham3/ASTGCN/tree/master/data>.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

## 13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We have included enough documentation in section 4, 5 and appendix B with respect to our introduced benchmarks.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.



- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

#### 14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

#### 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.