

1 **To Reviewer 1:** Thanks for your review. (I) Fair comparisons: the datasets in our evaluation come from ASTGCN and
2 STSGCN. Our results can be cross-checked to examine the performance (e.g., compare ref [11]). The code and datasets
3 are also provided for reproducing our results. For other baselines (e.g., DCRNN, STGCN), we fine-tuned them to
4 choose the best settings and hyper-parameters to ensure our comparison is fair. We believe our results are solid and the
5 superiority can be explained from two views: first, the predefined graph, which is provided by ASTGCN and STSGCN,
6 may be sub-optimal compared to our adaptively learned graph; second, the NAPL module can do unique fitting for each
7 node effectively. (II) Parameter size: we compared the number of parameters in Table 3 (page 8), which shows that
8 AGCRN has moderate parameters when compared with baselines (e.g., ASTGCN). Besides, if setting the embedding
9 dimension to a smaller value, AGCRN will contain much fewer parameters but still achieve good performance (as shown
10 in Figure 4). (III) E_g : there is no need to force E_g be orthogonal since an identity matrix is added to the GCN. (IV)
11 Learned graph: the values in the learned adaptive graph vary from 0.0014 to 0.0396 (for PeMSD4), which demonstrates
12 that DAGG can reveal relative importance among different series. We will visualize the learned adaptive graph with
13 heatmap and plot some highly-correlated series in the supplementary. We can achieve sparse graph and avoid negative
14 links by setting small values to zero. However, we didn't see it brings obvious improvements in our experiments. The
15 learned graph in DAGG can perform well and we will keep exploring more advanced graph generation methods.

16 **To Reviewer 2:** Thanks for your careful review. (I) Experiments: In Table 1, the "*" sign should be presented to
17 the results of STSGCN. As you have listed, DCRNN, STGCN, and ASTGCN used different datasets for evaluation.
18 We followed the most recent works (e.g., ASTGCN and STSGCN) and conducted experiments on the PeMSD4 and
19 PeMSD8 datasets. Our results can be cross-checked and compared with the results presented in ASTGCN and STSGCN.
20 (II) Metr-LA: the Metr-LA dataset is not the first choice for our evaluation because prior works in Metr-LA (e.g.,
21 DCRNN) use early fusion (e.g., concatenation) to integrate the extraneous data (e.g., time of day), which is not the
22 focus of our work and does not help the graph generation process. However, we still conduct experiments on Metr-LA
23 following your suggestion and compare our results with the reported results in DCRNN (will add to the supplementary).
24 Experiments show that AGCRN achieves comparable performance to DCRNN and consistently outperforms GCRNN
25 (i.e., DCRNN with undirected graph). The long-term prediction performance of AGCRN is especially superior (for 1
26 hour ahead prediction, MAE:3.57, RMSE:7.38, MAPE:10.09%). Our further analysis show AGCRN can gain more
27 advantages over DCRNN under a fair comparison (i.e., with the extraneous data removed for both models), even
28 though AGCRN does not rely on encoder-decoder, teacher forcing, or learning rate decay strategies while DCRNN
29 does. (III) Our preliminary experiments show that pair-wise distance or similarity-based graph can work together with
30 DAGG under a multi-graph framework and provide extra information. Also, we agree that generating outputs in a
31 sequential manner (e.g., under the encoder-decoder framework) may improve the results despite sacrificing efficiency.
32 It is convincing that AGCRN has potentials to enable more advanced performance in time series analysis. (IV) Indeed,
33 learning in large graphs for evolving systems is important yet not well studied in literature. Graph partition and
34 sub-graph training may probably work with some adaption. We will highlight this direction as a future work.

35 **To Reviewer 3:** Thanks for your comments. (I) Loss function: unfortunately, there is no standard about which loss
36 function is better. DCRNN uses L1 loss; STSGCN uses Huber Loss; STGCN uses L2 loss. We chose L1 loss because
37 it achieved better results in our experiments than L2 loss. We would like to highlight that we have carefully chosen
38 the best settings (e.g., loss function) for the baselines to make sure our comparison is fair unless stated otherwise. (II)
39 Prediction Length: first, each step is 5 mins in our experiments (see line 196 of page 5). Thus, we predict the future
40 60 mins data with the historical 60 mins data, which is the most common setting used for traffic forecasting. Second,
41 our problem is formulated to predict several future steps data with multiple historical steps data (Section 3.1). Thus,
42 the exact time lengths of the historical time period and future period are dependent on the number of time steps and
43 the length of one step, which could vary by applications. (III) Multi-step prediction: Yes, the results in Section 4.4
44 come directly from the results of multi-step prediction. As formulated in Section 3.1 and mentioned in Section 3.5, our
45 method can generate multi-step predictions directly. With 5 mins as the slot length in our experiments, the short-term (5
46 mins) and long-term (60 mins) predictions refer to the predictions made at the first step and the 12th step, respectively.

47 **To Reviewer 4:** Thanks for your comments. (I) we omitted some works about GCN for classification tasks due to the
48 space limitation and different problem settings. We will add a subsection under Section 2 to discuss those omitted
49 studies. Specifically, our work deal with dynamically evolving streams from nodes in the whole graph but GCN-based
50 classification models deal with static features/attributes from nodes within a neighborhood. Our work may share some
51 high-level similarities with these studies but has totally different designs for different purposes. For NAPL, our work
52 aim to learn unique parameters for each node, while GAT aims at learning different importance scores for neighbors.
53 The weights in GAT are still shared among different sub-graphs. The matrix factorization in NAPL may be similar with
54 clustering in graph pooling at a high level, but they work in different ways. For DAGG, the learned node embedding in
55 our work can be regarded as high-level abstraction of the whole data stream and is used to generate the graph. But in
56 GCN-based classification models, the learned node embedding is the transformation of nodes feature with the help of
57 existing graph. (II) we apologise for the typos. Metr-La should be PeMSD8 and STSGCN should be ref [11]. We have
58 revised them and conducted a thorough proofreading to eliminate possible typos and grammar mistakes.