

Lightweight and Interpretable Transformer via Mixed Graph Algorithm Unrolling for Traffic Forecast

Ji Qi¹ Mingxiao Liu¹ Tam Thuc Do² Yuzhe Li¹ Zhuoshi Pan¹
Gene Cheung^{2*} H. Vicky Zhao^{1*}

¹Department of Automation, Tsinghua University, Beijing, China

²Department of EECS, York University, Toronto, Canada

1 June, 2026



Challenges in Traffic forecasting: modeling complex spatiotemporal dependencies, limited training data, edge device deployment, etc.

Conventional Transformers

- Data-driven approach with complex dependencies modeled by attention
- Powerful performance
- Large model size and computational cost, lack of interpretability

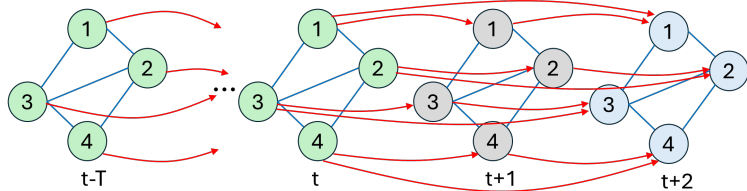
Graph Signal Processing (GSP)

- Model-based approach by smoothing signals on graphs
- Strong interpretability
- Limited performance and flexibility

Motivation: Building Lightweight and Interpretable Transformers

- *Parameterize* the GSP-based approach for data-driven learning while preserving the interpretability — **Algorithm Unrolling**
- Preserve the *transformer-like architecture* by periodically introducing graph learning modules to learn graphs from data

Mixed-Graph Modeling & Smoothing Priors



- **Directed graphs** \mathcal{G}^d : temporal dependencies on each node with its previous W states
- **Undirected graphs** \mathcal{G}^u : spatial dependencies between nodes with road connections on each timestamp

Mixed Graph Smoothing Priors

- Undirected graph: $\text{GLR} = \mathbf{x}^\top \mathbf{L}^u \mathbf{x}$
- Directed graph: $\text{DGLR} = \mathbf{x}^\top (\mathbf{L}_r^d)^\top \mathbf{L}_r^d \mathbf{x} \triangleq \mathbf{x}^\top \mathcal{L}_r^d \mathbf{x}$ and $\text{DGTV} = \|\mathbf{L}_r^d \mathbf{x}\|_1$

Smoothing the signals on the mixed graph:

$$\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_2^2 + \mu_u \mathbf{x}^\top \mathbf{L}^u \mathbf{x} + \mu_{d,2} \mathbf{x}^\top \mathcal{L}_r^d \mathbf{x} + \mu_{d,1} \|\mathbf{L}_r^d \mathbf{x}\|_1. \quad (1)$$

Introduce the auxiliary variables: $\phi = \mathbf{L}_r^d \mathbf{x}$ for the ℓ_1 term, **split the ℓ_2 terms** and add **auxiliary variables** $\mathbf{z}_u, \mathbf{z}_d$

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{z}_u, \mathbf{z}_d} \quad & \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_2^2 + \mu_u \mathbf{z}_u^\top \mathbf{L}^u \mathbf{z}_u + \mu_{d,2} \mathbf{z}_d^\top \mathcal{L}_r^d \mathbf{z}_d + (\gamma^\tau)^\top (\phi^\tau - \mathbf{L}_r^d \mathbf{x}) + \frac{\rho}{2} \|\phi^\tau - \mathbf{L}_r^d \mathbf{x}\|_2^2 \\ \text{s.t.} \quad & \mathbf{x} = \mathbf{z}_u = \mathbf{z}_d, \quad \phi = \mathbf{L}_r^d \mathbf{x}. \end{aligned} \quad (2)$$

Write the augmented Lagrangian function, and solve with ADMM:

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{z}_u, \mathbf{z}_d, \phi} \quad & \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_2^2 + \mu_u \mathbf{z}_u^\top \mathbf{L}^u \mathbf{z}_u + \mu_{d,2} \mathbf{z}_d^\top \mathcal{L}_r^d \mathbf{z}_d + (\gamma^\tau)^\top (\phi^\tau - \mathbf{L}_r^d \mathbf{x}) + \frac{\rho}{2} \|\phi^\tau - \mathbf{L}_r^d \mathbf{x}\|_2^2 \\ & + (\gamma_u^\tau)^\top (\mathbf{x} - \mathbf{z}_u) + \frac{\rho_u}{2} \|\mathbf{x} - \mathbf{z}_u\|_2^2 + (\gamma_d^\tau)^\top (\mathbf{x} - \mathbf{z}_d) + \frac{\rho_d}{2} \|\mathbf{x} - \mathbf{z}_d\|_2^2 + \mu_{d,1} \|\phi^\tau\|_1 \end{aligned} \quad (3)$$

① **Optimize $\mathbf{x}^{\tau+1}$ and auxiliary variables $\mathbf{z}_u^{\tau+1}, \mathbf{z}_d^{\tau+1}$:**

Solve the following linear systems by conjugated gradient (CG) descent:

$$(\mathbf{H}^\top \mathbf{H} + \frac{\rho}{2} \mathcal{L}_r^d + \frac{\rho_u + \rho_d}{2} \mathbf{I}) \mathbf{x}^{\tau+1} = (\mathbf{L}_r^d)^\top \left(\frac{\gamma^\tau}{2} + \frac{\rho}{2} \phi^\tau \right) - \frac{\gamma_u^\tau}{2} + \frac{\rho_u}{2} \mathbf{z}_u^\tau - \frac{\gamma_d^\tau}{2} + \frac{\rho_d}{2} \mathbf{z}_d^\tau + \mathbf{H}^\top \mathbf{y}, \quad (4)$$

$$(\mu_u \mathbf{L}^u + \frac{\rho_u}{2} \mathbf{I}) \mathbf{z}_u^{\tau+1} = \frac{\gamma_u^\tau}{2} + \frac{\rho_u}{2} \mathbf{x}^{\tau+1}, \quad (5)$$

$$(\mu_{d,2} \mathcal{L}_r^d + \frac{\rho_d}{2} \mathbf{I}) \mathbf{z}_d^{\tau+1} = \frac{\gamma_d^\tau}{2} + \frac{\rho_d}{2} \mathbf{x}^{\tau+1}. \quad (6)$$

Structural Interpretation: Alternative *low-pass filtering* with undirected graph kernel \mathbf{L}^d and directed graph kernels $\mathbf{L}_r^d, \mathcal{L}_r^d$

$$\begin{aligned} \mathbf{x}^{\tau+1} &= \mathbf{U} \text{diag} \left(\left\{ \frac{2}{\rho \xi_k + \rho_u + \rho_d + 2 \cdot \mathcal{K}_{\text{obs}}} \right\} \right) \mathbf{U}^\top \text{RHS}_{(4)}, \\ \mathbf{z}_u^{\tau+1} &= \mathbf{V} \text{diag} \left(\left\{ \frac{2}{2\mu_u \lambda_k + \rho_u} \right\} \right) \mathbf{V}^\top \left(\frac{\gamma_u^\tau}{2} + \frac{\rho_u}{2} \mathbf{x}^{\tau+1} \right), \\ \mathbf{z}_d^{\tau+1} &= \mathbf{U} \text{diag} \left(\left\{ \frac{2}{2\mu_{d,2} \xi_k + \rho_d} \right\} \right) \mathbf{U}^\top \left(\frac{\gamma_d^\tau}{2} + \frac{\rho_d}{2} \mathbf{x}^{\tau+1} \right) \end{aligned}$$

- ② **Optimize $\phi^{(t+1)}$:** The $\ell_2 - \ell_1$ problem gives a closed-form solution, where S is the soft-thresholding operator:

$$\delta = (\mathbf{L}_r^d)_i \mathbf{x}^{\tau+1} - \frac{1}{\rho} \gamma_i^\tau, \quad \phi_i^{\tau+1} = \text{sign}(\delta) \cdot \max(|\delta| - \rho^{-1} \mu_{d,1}, 0) \triangleq S_{\mu_{d,1}/\rho}(\delta). \quad (7)$$

Structural Interpretation: the high-frequency components $\mathbf{L}_r^d \mathbf{x}$ of signal \mathbf{x} is *attenuated* by the soft-thresholding operation — *low-pass filtering* of signal \mathbf{x} with the directed graph kernel \mathbf{L}_r^d under *filter-bank explanation*

- ③ **Lagrangian Multipliers $\gamma, \gamma_u, \gamma_d$ update:**

$$\begin{aligned} \gamma^{\tau+1} &= \gamma^\tau + \rho(\phi^{\tau+1} - \mathbf{L}_r^d \mathbf{x}^{\tau+1}), \\ \gamma_u^{\tau+1} &= \gamma_u^\tau + \rho_u(\mathbf{x}^{\tau+1} - \mathbf{z}_u^{\tau+1}), \\ \gamma_d^{\tau+1} &= \gamma_d^\tau + \rho_d(\mathbf{x}^{\tau+1} - \mathbf{z}_d^{\tau+1}) \end{aligned} \quad (8)$$

Unrolling the ADMM algorithm yields an ADMM block interpreted as **a sequence of learnable low-pass filters** over the mixed graph \implies **FFN** in transformer.

Conventional Attention Mechanism:

$$e(i, j) = (\mathbf{Q}\mathbf{x}_j)^\top (\mathbf{K}\mathbf{x}_i) = \mathbf{x}_j^\top (\mathbf{Q}^\top \mathbf{K}) \mathbf{x}_i, \quad a_{i,j} = \frac{\exp(e(i,j))}{\sum_{l=1}^N \exp(e(i,l))}, \quad \mathbf{y}_i = \sum_{l=1}^N a_{i,l} \mathbf{x}_l \mathbf{V} \quad (9)$$

Lightweight Graph Learning Module

- **Feature extraction:** leveraging neighboring node signals and position information, with learnable function $f(\cdot)$:

$$\mathbf{f}_i = f(i | \mathbf{X}_{\mathcal{N}_i}, \mathbf{E}_{\mathcal{N}_i}) \quad (10)$$

- **Graph construction:** Gaussian kernel under Mahalanobis distances, with learnable matrices \mathbf{M} , \mathbf{P} :

$$d^u(i, j) = (\mathbf{f}_i^u - \mathbf{f}_j^u)^\top \mathbf{M} (\mathbf{f}_i^u - \mathbf{f}_j^u), \quad d^d(i, j) = (\mathbf{f}_i^d - \mathbf{f}_j^d)^\top \mathbf{P} (\mathbf{f}_i^d - \mathbf{f}_j^d) \quad (11)$$

- **Normalization:**

$$w_{i,j}^u = \frac{\exp(-d^u(i,j))}{\sqrt{\sum_{l \in \mathcal{N}_i} \exp(-d^u(i,l))} \sqrt{\sum_{k \in \mathcal{N}_j} \exp(-d^u(k,j))}}, \quad w_{j,i}^d = \frac{\exp(-d^d(i,j))}{\sum_{[k,j] \in \mathcal{E}^d} \exp(-d^d(k,j))}. \quad (12)$$

The Unrolled Network Architecture

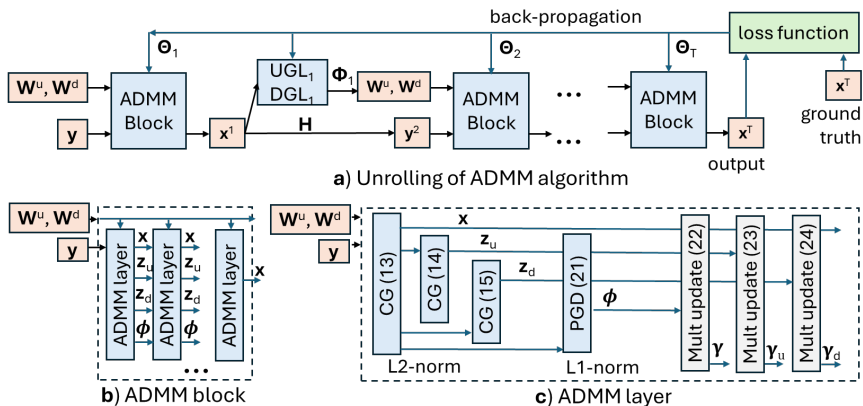


Figure 2: Unrolling of proposed iterative ADMM algorithm into neural layers.

Loss function: Huber loss on the *entire* output sequence $\hat{\mathbf{X}}_{t-T:t+S}$ for the **signal recovery** problem.

Experiment Setup & Results

- **Tasks:** 30 / 60 / 120-minute traffic forecasting on real world datasets: PEMS03 (358 nodes, 547 edges) and METR-LA (207 nodes, 1,315 edges)
- **Evaluation metrics:** MAE, RMSE, MAPE (%)
- **Baselines (with parameter count comparison) :**

Category	Selected Models	Params #
Model-based	VAR (2003)	-
GNNs	STGCN (2018), STSGCN (2020)	321K, 3,496K
GATs	GMAN (2020), ST-Wave (2023)	210K, 883K
Transformers	PDFormer (2023), STAEformer (2023), PatchSTG (2025)	531K, 1,404K, 2,283K
Adaptive GNN	Graph WaveNet (2019), AGCRN (2020)	277K, 749K
MLP-based models	STID (2022), SimpleTM (2025)	123K, 540K
Mixed-graph unrolling	Ours	38K

Our model employs drastically fewer parameters than all baselines (7.2% of *PDFormer*).

Experimental Results

Our model achieves *competitive performance* (top 3 in 4 out of 6 runs) with *significantly fewer parameters*.

Table 1: Comparison of metrics of our lightweight transformer to the baselines. We use **boldface** for the smallest error, color the 2nd and 3rd small error in **blue**, and underline the next 2 smallest errors.

Dataset & Horizon	PEMS03 (358 nodes, 547 edges)			METR-LA (207 nodes, 1,315 edges)		
	30 minutes	60 minutes	120 minutes	30 minutes	60 minutes	120 minutes
VAR	28.07 / 16.53 / 17.49	30.54 / 18.31 / 19.61	36.65 / 22.40 / 24.64	10.72 / 5.55 / 11.29	12.59 / 6.99 / 13.46	14.83 / 8.90 / 16.54
STGCN	28.06 / 18.24 / 17.76	37.31 / 24.29 / 23.00	54.34 / 34.83 / 30.52	11.26 / 5.08 / 10.93	13.91 / 6.35 / 13.67	16.92 / 8.11 / 17.69
STSGCN	26.41 / 16.75 / 16.86	30.62 / 19.35 / 19.15	38.04 / 23.86 / 23.62	10.25 / 4.05 / 9.18	12.65 / <u>5.18</u> / 11.48	15.74 / 6.84 / 15.33
GMAN	25.79 / 16.23 / 20.04	<u>27.57</u> / 17.48 / 24.33	30.69 / <u>19.20</u> / 26.69	11.97 / 5.34 / 10.66	14.49 / 6.93 / 13.12	15.53 / 7.59 / 14.70
ST-Wave	<u>25.57</u> / 15.11 / 15.04	28.65 / <u>16.81</u> / 19.24	29.88 / 17.11 / 17.71	10.81 / 4.11 / 9.16	13.24 / 5.33 / <u>11.32</u>	23.18 / 11.22 / 12.71
PDFormer	23.71 / 15.05 / 18.16	<u>27.16</u> / 17.26 / 21.21	35.77 / 22.25 / 25.01	<u>10.21</u> / 3.89 / 8.50	12.17 / 4.81 / 10.92	17.27 / 9.55 / 19.10
STAEformer	30.22 / 18.85 / 26.62	38.36 / 23.68 / 29.21	48.72 / 31.96 / 44.71	<u>10.16</u> / 3.73 / 8.25	12.58 / 4.79 / 10.08	14.63 / 5.82 / 11.87
PatchSTG	22.89 / 14.66 / 17.51	25.09 / 15.76 / 17.02	29.97 / 18.65 / 21.13	10.24 / 3.89 / <u>9.02</u>	14.12 / 5.80 / 12.08	17.39 / 7.73 / 16.48
GraphWaveNet	25.76 / 15.22 / <u>16.83</u>	28.15 / 18.11 / <u>17.56</u>	34.60 / 21.10 / <u>22.45</u>	11.42 / 5.18 / 8.21	<u>12.46</u> / <u>5.17</u> / 11.97	23.13 / 11.32 / <u>13.53</u>
AGCRN	27.40 / <u>15.19</u> / 14.35	29.90 / <u>16.76</u> / 15.32	<u>32.79</u> / 18.72 / 17.03	10.11 / 3.82 / 8.61	<u>12.56</u> / 5.00 / 11.25	14.77 / 6.33 / 14.05
STID	26.50 / 17.27 / 18.59	31.88 / 20.82 / 24.89	42.98 / 28.03 / 42.41	<u>10.21</u> / <u>4.05</u> / 9.47	12.61 / 5.21 / 12.22	15.48 / <u>6.98</u> / 16.78
SimpleTM	23.75 / 15.13 / 15.59	<u>25.56</u> / 15.97 / 15.49	30.58 / 18.73 / 18.18	8.76 / 4.12 / 7.63	11.96 / 5.49 / 9.65	<u>15.44</u> / 7.64 / 12.27
Ours	<u>25.05</u> / 15.85 / <u>16.49</u>	26.96 / 16.58 / <u>18.07</u>	34.06 / 20.23 / 23.86	10.06 / 4.05 / 9.23	12.17 / 5.27 / 11.78	<u>15.19</u> / <u>7.14</u> / 16.44

Experimental Results

Our model sits distinctly in the lower-left corner in the error-size trade-off plot, achieving the accuracy of models over $10\times$ larger.

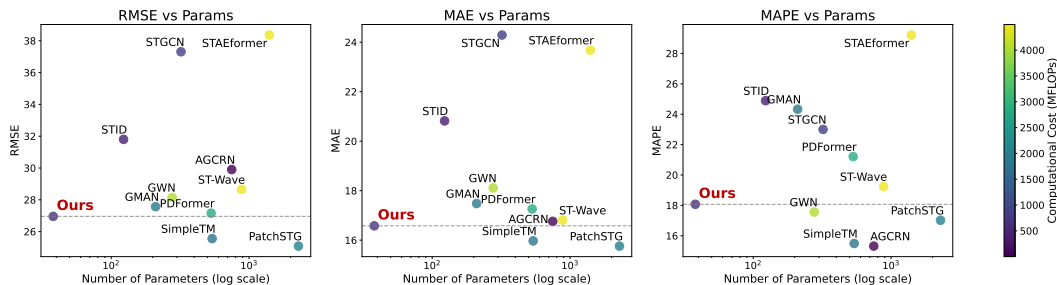


Figure 3: The trade-off between the number of parameters and performance for 60-minute forecast on PEMS03 dataset.

- We unroll a mixed-graph smoothing algorithm into a lightweight and interpretable transformer architecture for traffic forecasting, achieving competitive performance with significantly fewer parameters.
- **Broader impacts:**
 - Transformers can be *whitened* and *reduced* as an unrolled network of graph smoothing algorithms.
 - We can learn *directed* attention scores with directed graph modeling, which enhance the generalizability by capturing assymetrical or causal dependencies within data.

Thanks!