

WAVE: Weighted Autoregressive Varying Gate Attention for Time-Series Forecasting



Jiecheng Lu¹ Xu Han² Yan Sun¹ Shihao Yang¹

¹Georgia Institute of Technology ²Amazon Web Services

Motivation

- Decoder-only Transformers (GPT-style) are powerful for sequential modelling but rarely used *end-to-end* in long-horizon time-series forecasting (TSF).
- Classic ARMA models decouple long-term trends (AR) from short-term fluctuations (MA).
- Existing efficient attentions add exponential decay (EMA) for locality, but this can suppress stable seasonal patterns.

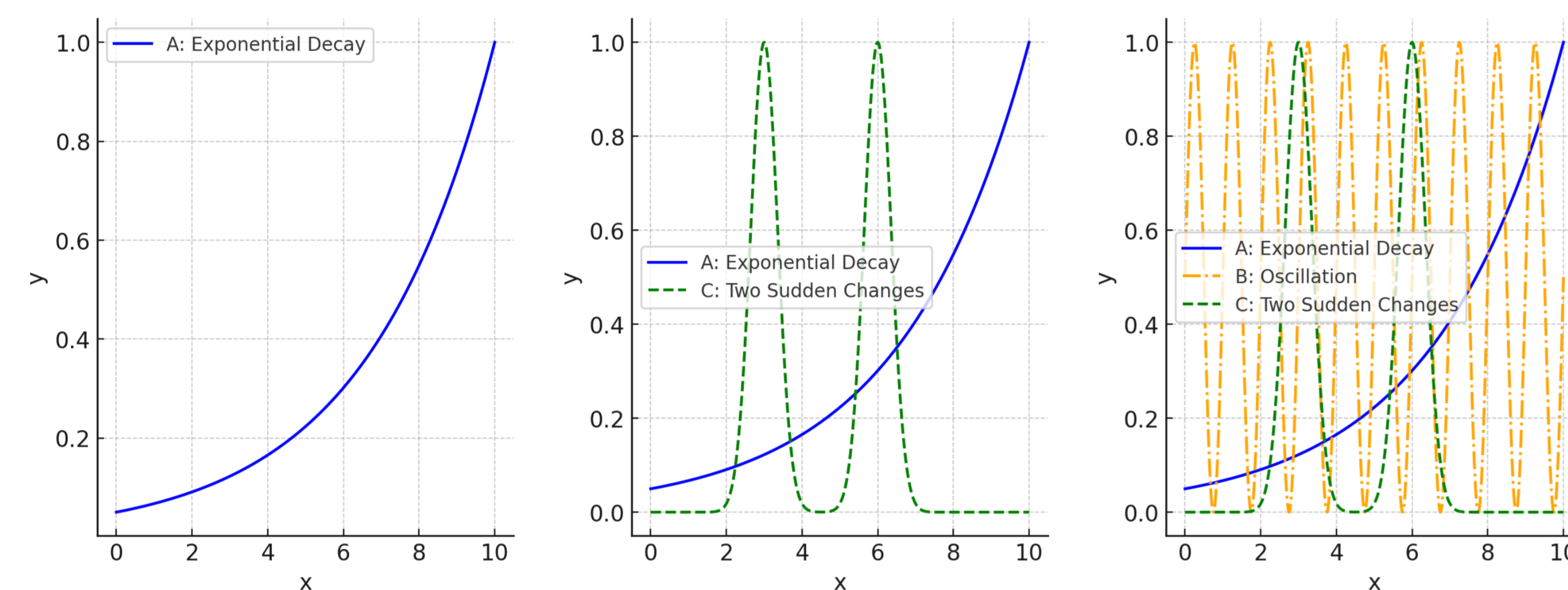


Figure 1. Visualization of different effects with exponential decay strategies and their challenges in gated linear attention. (Left: a) Pure exponential decay strategy in gated linear attention; (Mid: b) Exponential decay facing challenges in capturing long-term dependencies; (Right: c) Exponential decay facing challenges in capturing periodic dependencies

Key Contributions

- Show that a well-tokenised AR Transformer already reaches SOTA TSF performance.
- Propose **WAVE attention**: adds an MA path to any AR attention (softmax, linear, gated linear, ...).
- Novel *indirect MA weight generation* keeps $\mathcal{O}(N)$ time and no extra parameters (weights are shared).
- Achieve new state-of-the-art results on 12 public TSF benchmarks.

One-Step Patch Tokenization

Each multivariate series is split into non-overlapping *patch tokens* of size L_P . A single decoder step predicts the next patch, preventing error accumulation.

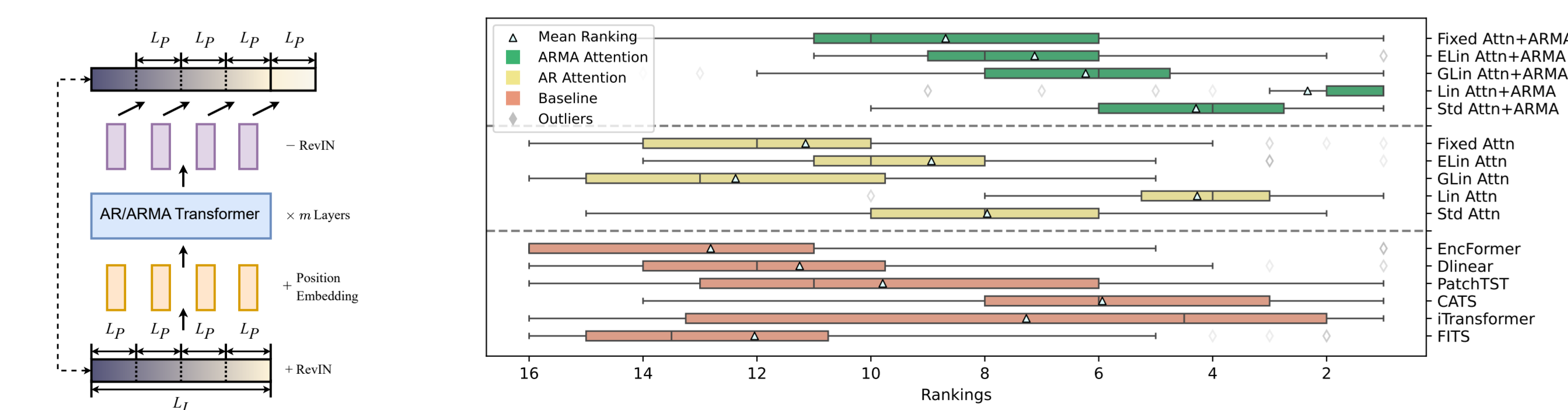


Figure 2. (Left: a) Overall architecture of our decoder Transformer for TSF. (Right: b) Box plots of performance rankings from 48 sub-experiments across 12 datasets. Green represents WAVE Transformers, yellow AR Transformers, and red the baselines, with triangles indicating mean rankings. AR Transformers perform comparably to baselines, while WAVE Transformers significantly outperform their AR counterparts. See Table and for more details.

WAVE Attention Mechanism

Given query \mathbf{q}_t , keys $\mathbf{k}_{1:t}$, values $\mathbf{v}_{1:t}$:

$$\mathbf{o}_t^{\text{AR}} = \sum_{j=1}^t w_{t,i} \mathbf{v}_j \quad + \quad \mathbf{o}_t^{\text{MA}} = \sum_{j=1}^{t-1} \theta_{t-1,j} \boldsymbol{\epsilon}_j$$

AR (long-term) MA (short-term)

$$\boldsymbol{\epsilon}_j = \mathbf{v}_{j+1} - \mathbf{o}_j^{\text{AR}}$$

The MA weights $\theta_{t-1,j}$ are produced *indirectly*: $\theta = B(I-B)^{-1}$, where $B_{t-1,j} = \phi_q(\mathbf{q}_{t-1}^{\text{MA}}) \phi_k(\mathbf{k}_j^{\text{MA}})$.

- AR path**: standard causal attention (softmax / linear / gated).
- MA path**: lightweight linear attention over residuals $\boldsymbol{\epsilon}$.
- Weight-sharing**: \mathbf{W}_q reused; \mathbf{W}_v fixed to identity.

Indirect MA Weight Generation

(1) Avoids explicit $N \times N$ MA matrix. (2) Preserves $\mathcal{O}(N)$ runtime of efficient attentions. (3) Activation choice: $\phi_k = \sigma(\alpha k)$, $\phi_q = -\text{LeakyReLU}(-q) \rightarrow$ yields near-diagonal Θ emphasising local context.

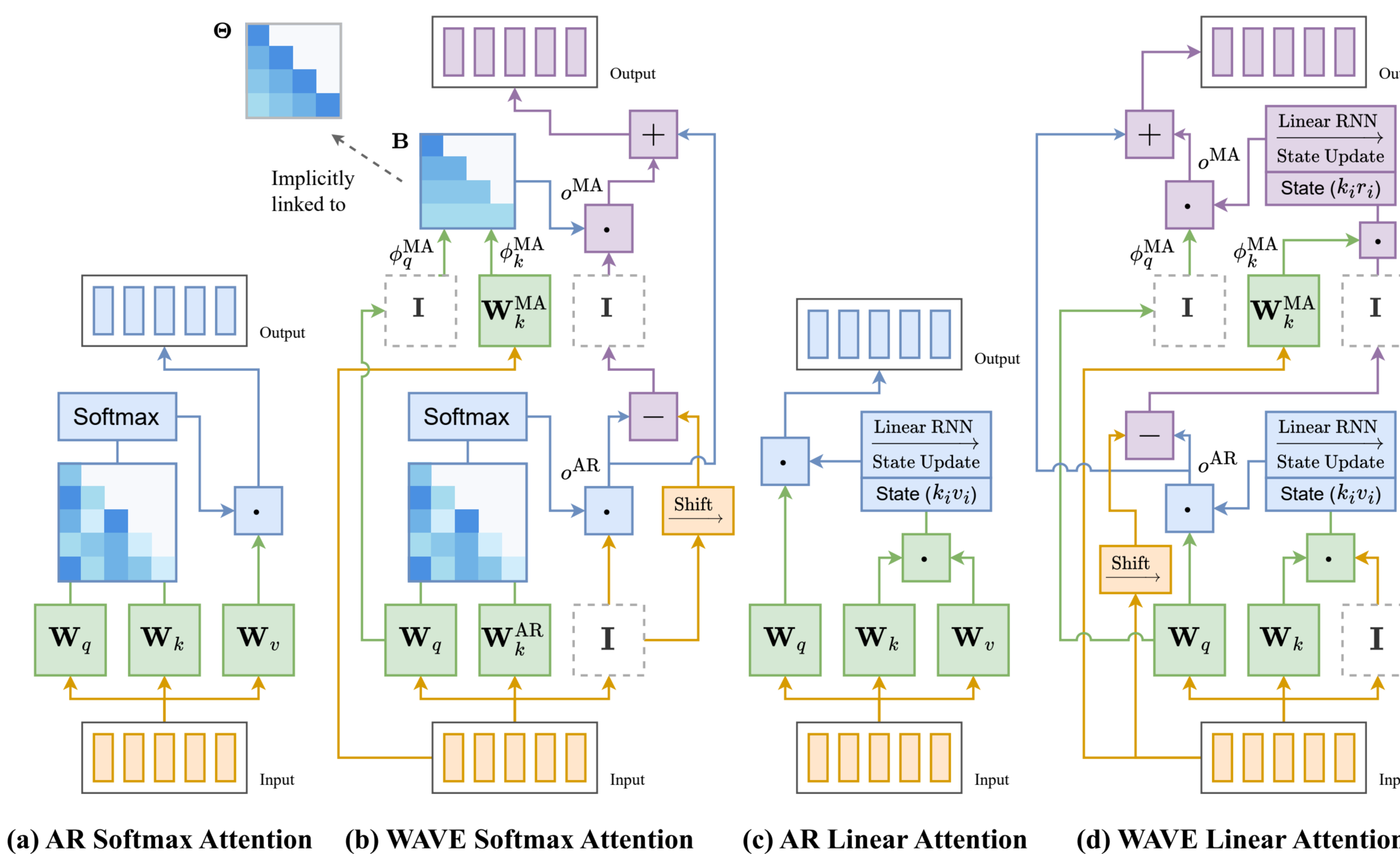


Figure 3. WAVE attention structure with the indirect MA weight generation method applied to softmax and linear attention. See Table 1 for more calculation details.

Table 1. Summary of WAVE attention for various attention mechanisms, detailing the calculation methods for AR output and MA output, where $\mathbf{r}_j = \mathbf{v}_{j+1} - \mathbf{o}_j^{\text{AR}}$.

Model	AR term output \mathbf{o}_t^{AR}	Indirect MA term output \mathbf{o}_t^{MA}
Standard Softmax Attention (Std Attn)	$\sum_{i=1}^t \frac{\exp(\mathbf{q}_t(\mathbf{k}_i^{\text{AR}})^T) \mathbf{v}_i}{\sum_{i=1}^t \exp(\mathbf{q}_t(\mathbf{k}_i^{\text{AR}})^T)}$	$\sum_{j=1}^{t-1} \phi_q^{\text{MA}}(\mathbf{q}_{t-1}) \phi_k^{\text{MA}}(\mathbf{k}_j^{\text{MA}})^T \mathbf{r}_j$
Linear Attention (Lin Attn)	$\mathbf{q}_t \sum_{i=1}^t (\mathbf{k}_i^{\text{AR}})^T \mathbf{v}_i$	$\phi_q^{\text{MA}}(\mathbf{q}_{t-1}) \sum_{j=1}^{t-1} \phi_k^{\text{MA}}(\mathbf{k}_j^{\text{MA}})^T \mathbf{r}_j$
Element-wise Linear Attention (ELin Attn)	$\sigma(\mathbf{q}_t) \odot \sum_{i=1}^t \frac{\exp(\mathbf{k}_i^{\text{AR}} \odot \mathbf{v}_i)}{\sum_{i=1}^t \exp(\mathbf{k}_i^{\text{AR}})}$	$\phi_q^{\text{MA}}(\mathbf{q}_{t-1}) \odot \sum_{j=1}^{t-1} \phi_k^{\text{MA}}(\mathbf{k}_j^{\text{MA}}) \odot \mathbf{r}_j$
Gated Linear Attention (GLin Attn)	$\mathbf{q}_t \sum_{i=1}^t \mathbf{G}_i \odot (\mathbf{k}_i^{\text{AR}})^T \mathbf{v}_i$	$\phi_q^{\text{MA}}(\mathbf{q}_{t-1}) \sum_{j=1}^{t-1} \phi_k^{\text{MA}}(\mathbf{k}_j^{\text{MA}})^T \mathbf{r}_j$
Fixed Attention (Fixed Attn)	$\sum_{i=1}^t w_{t,i}^{\text{AR}} \mathbf{v}_i$	$\phi_q^{\text{MA}}(\mathbf{w}_{t-1}^{\text{MA},q}) \sum_{j=1}^{t-1} \phi_k^{\text{MA}}(\mathbf{w}_j^{\text{MA},k})^T \mathbf{r}_j$

Main Results (12 Datasets)

- Average ranking over 48 TSF tasks** (4 horizons \times 12 data sets): Δ AR Transformer rank = 4.3 \star **WAVE-Linear rank = 2.3 (best)**
- WAVE improves every attention variant: e.g. on “Weather”, MSE 4% vs AR; on “ETTm1”, MSE 8%.
- Outperforms previous models like PatchTST, FITS, iTransformer, CATS, DLinear.

Table 2. Summary of main TSF results with forecasting horizons $L_P \in \{12, 24, 48, 96\}$ and $L_I = 512$. Average rankings (AvgRank) of each model on test set MSE, along with the count of first-place rankings (#Top1), are included.

Models	Std Attn	Std Attn +ARMA	Lin Attn	Lin Attn +ARMA	GLin Attn	GLin Attn +ARMA	ELin Attn	ELin Attn +ARMA	Fixed Attn	Fixed Attn +ARMA
AvgRank	7.958	<u>4.292</u>	4.271	<u>2.333</u>	12.375	<u>6.229</u>	8.938	<u>7.125</u>	11.146	<u>8.688</u>
#Top1	0	<u>4</u>	4	<u>25</u>	0	<u>1</u>	1	<u>3</u>	1	<u>3</u>

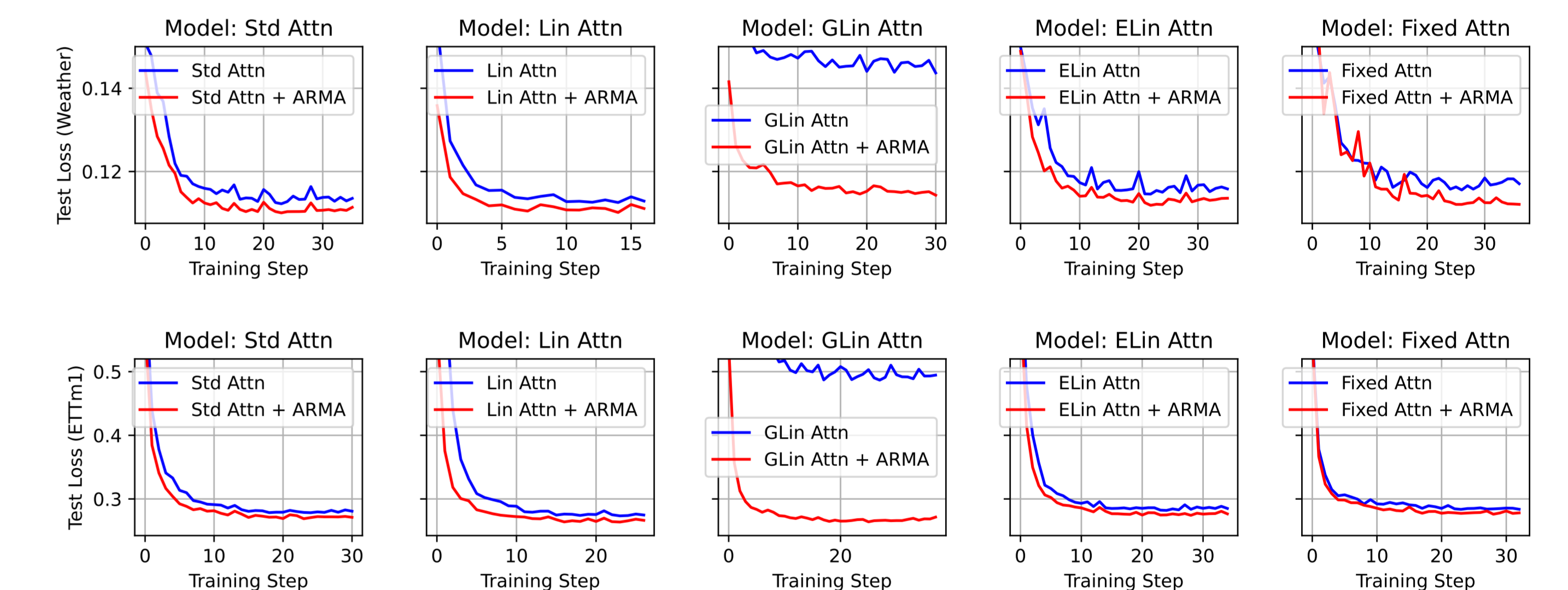


Figure 4. Visualization of test loss curves. We show the testing performance of five attention mechanisms using WAVE structures on the Weather and ETTm1 datasets ($L_I = 512$, $L_P = 48$).

Scalability & Efficiency

- Longer look-back ($L_I=4096$) \rightarrow WAVE MSE improves; many baselines degrade.
- WAVE with 3 layers beats AR with 8 layers \rightarrow gains come from *structure*, not size.
- Same parameter count, linear memory/time; negligible extra FLOPs.

Conclusion

WAVE bridges classic ARMA ideas and modern Transformers. By explicitly modeling residuals via an efficient MA branch, it decouples short- and long-term dynamics, giving state-of-the-art forecasting accuracy without extra complexity.