# Memformer

## A Memory Guided Transformer for Time Series Forecasting

Yunyao Cheng, Chenjuan Guo, Bin Yang, Haomin Yu, Kai Zhao, Christian S. Jensen

February 2025

*Presented by* **Andreas Gottschalk Krath**

# 1. Introduction

**Forecasting**

- Predicting the future
  - ‣ Allows preparation
- Many applications
  - ‣ Electricity prices
  - ‣ Finance
- Long term forecasting?
  - ‣ Obviously more difficult than short term
  - ‣ Time constrained tasks

**Long Term Forecasting**

- What defines long term?
  - ▸ Historical horizon
  - ▸ Forecasting horizon
  - ▸ Both exceed 96 time steps
    - – Hourly time step $\rightarrow$ 4 days
  - ▸ Time series

**Variable Correlation**

- Complex systems have many variables
- $A$ increases and $B$ increases $\rightarrow$ Positive
- $A$ increases and $B$ decreases $\rightarrow$ Negative
- $A$ increases and $B$ is stagnant $\rightarrow$ None
- These impact forecasting accuracy
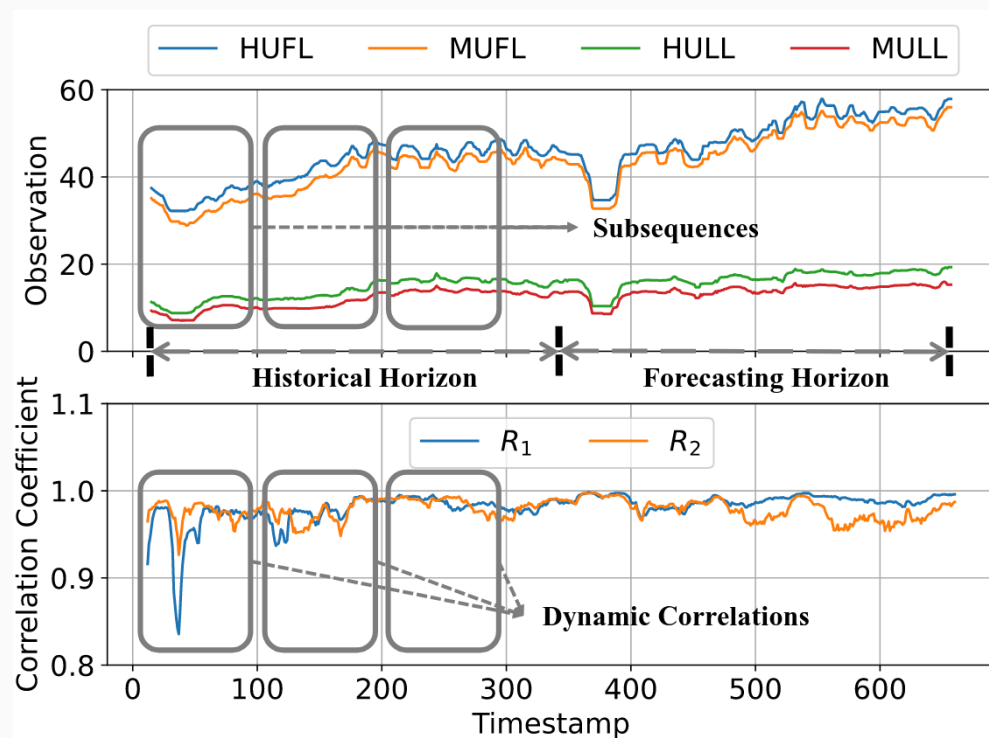  - ▸ Patterns in the data

**Dynamic Correlations**

- Are variable correlations stable over time?
  - ‣ No
- Correlations are dynamic over time
  - ‣ Seasons
  - ‣ Sensor drift
- We often consider average
  - ‣ Especially hurtful in time series
  - ‣ Predictions are bad in periods

## Dynamic Correlations

- Are variable correlations stable over time?
  - ▸ No
- Correlations are dynamic over time
  - ▸ Seasons
  - ▸ Sensor drift
- We often consider average
  - ▸ Especially hurtful in time series
  - ▸ Predictions are bad in periods



(a) Dynamic correlations. The Average $R_1 = 0.995$ and $R_2 = 0.990$.

**Disrupted Correlations**

- System errors
- External influence
- What happens with outliers?
  - ‣ Affect correlation $\longrightarrow$ accuracy
- Many models are sensitive to outliers
  - ‣ Numeric difference dominates training
  - ‣ Reason for a lot of preprocessing
    - – Normalization
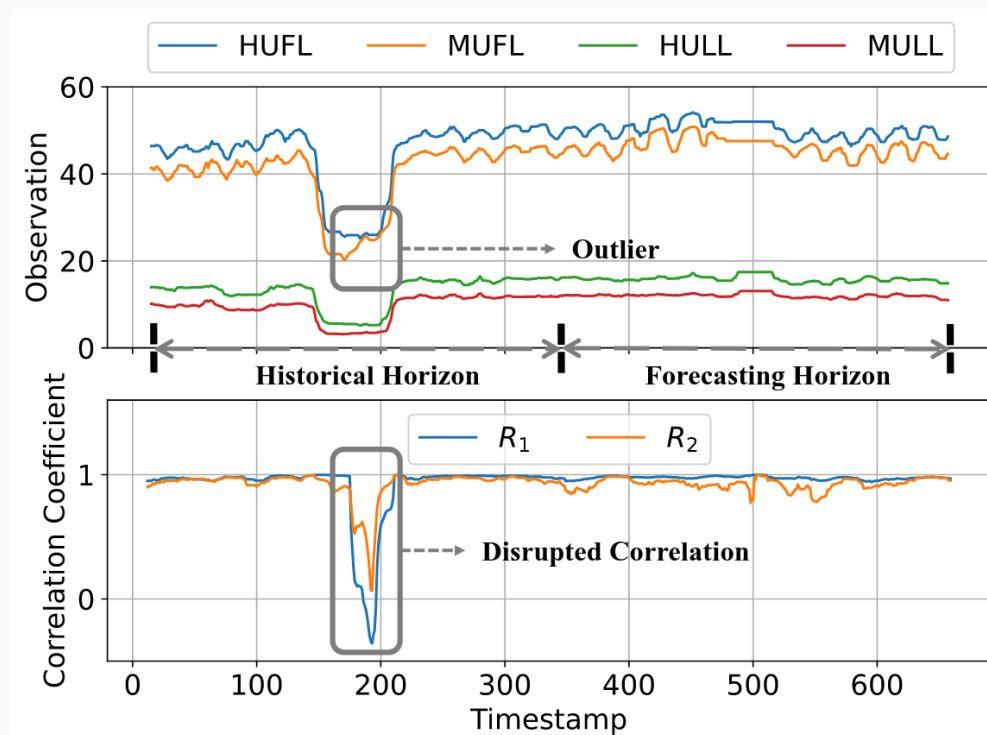    - – Clipping
    - – Pruning

## Disrupted Correlations

- System errors
- External influence
- What happens with outliers?
  - ▸ Affect correlation $\rightarrow$ accuracy
- Many models are sensitive to outliers
  - ▸ Numeric difference dominates training
  - ▸ Reason for a lot of preprocessing
    - – Normalization
    - – Clipping
    - – Pruning



(b) Disrupted correlation. The Average $R_1 = 0.908$ and $R_2 = 0.963$.

**Challenge 1**

- Capture dynamic correlations
- Mitigate disrupted correlations
- Existing solutions struggle with the latter
  - ▸ Capture dynamic and disrupted
  - ▸ Reduces model robustness

**Challenge 2**

- Local information 🤝 global information
- Global information is *all* local information
- Local information *affects* global information
- Existing solutions struggle with combining
  - ▸ Only local
  - ▸ Only global

**Memformer**

- Transformer
- Patch-wise recurrent graph learning
  - ▸ Captures dynamic correlations
- Global attention
  - ▸ Mitigates disrupted correlations
- Adresses challenge 1

**Alternating Memory Enhancer**

- Memory network
- Associates local and global information
- Adresses challenge 2

**Experiments**

- Proof

# 2. Methodology

**Instance normalization**

- Normalize within historical horizon only
- Mitigates the issue of internal covariate shift
- Allows model to effectively grasp the intricate temporal dynamics inherent in time series

$H' = (H - \mu)/\sqrt{(\sigma^2 + c)}, \text{where}$

$H$ is the historical horizon

$\mu$ is the mean

$\sigma$ is the variance

$c$ ensures numerical stability

**Instance normalization**

- Normalize within historical horizon only
- Mitigates the issue of internal covariate shift
- Allows model to effectively grasp the intricate temporal dynamics inherent in time series

$$H' = (H - \mu)/\sqrt{(\sigma^2 + c)}, \text{where}$$

$H$ is the historical horizon

$\mu$ is the mean

$\sigma$ is the variance

$c$ ensures numerical stability

poral dynamics inherent in time series. Instance normalization is defined as $\mathbf{H'} = (\mathbf{H} - \mu)/\sqrt{(\sigma^2 + \text{constant})}$, where $\mathbf{H'}$ denotes the preprocessed feature, $\mu$ and $\sigma$ denote the mean and variance of the sample, respectively, and "constant" is a small positive real number included to ensure numerical stability.

**Instance normalization**

- Normalize within historical horizon only
- Mitigates the issue of internal covariate shift
- Allows model to effectively grasp the intricate temporal dynamics inherent in time series

$H' = (H - \mu)/\sqrt{(\sigma^2 + c)}$, where

$H$ is the historical horizon

$\mu$ is the mean

$\sigma$ is the variance

$c$ ensures numerical stability

poral dynamics inherent in time series. Instance normalization is defined as $\mathbf{H'} = (\mathbf{H} - \mu)/\sqrt{(\sigma^2 + \text{constant})}$, where $\mathbf{H'}$ denotes the preprocessed feature, $\mu$ and $\sigma$ denote the mean and variance of the sample, respectively, and "constant" is a small positive real number included to ensure numerical stability.

- Mistake in variance?
  - ▸ $\sigma$ is conventional notation for standard deviation
  - ▸ $\sigma^2$ is conventional notation for variance

**What is going on?**

**What is going on?**

- Explored code to find answer
- `data_provider/data_loader.py`
  - ‣ Only place anything related to loading data happens
  - ‣ `Dataset_ETT_hour`, `Dataset_ETT_minute`, `Dataset_Custom`, `Dataset_Pred`

**What is going on?**

- Explored code to find answer
- `data_provider/data_loader.py`
  - ▸ Only place anything related to loading data happens
  - ▸ `Dataset_ETT_hour`, `Dataset_ETT_minute`, `Dataset_Custom`, `Dataset_Pred`

```python
from sklearn.preprocessing import StandardScaler
class ...:
    def __read_data__(self):
        self.scalar = StandardScaler()
        self.scaler.fit(train_data.values)
        data = self.scaler.transform(df_data.values)
```

**What is going on?**

- Explored code to find answer
- `data_provider/data_loader.py`
  - ‣ Only place anything related to loading data happens
  - ‣ `Dataset_ETT_hour`, `Dataset_ETT_minute`, `Dataset_Custom`, `Dataset_Pred`

```python
from sklearn.preprocessing import StandardScaler
class ...:
    def __read_data__(self):
        self.scalar = StandardScaler()
        self.scaler.fit(train_data.values)
        data = self.scaler.transform(df_data.values)
```

- They fit on training data
- Normalize entire dataset with μ and σ from training data

**What are they actually doing?**

<table>
<tr><td align="center">Preprocessing</td><td align="center">StandardScaler</td></tr>
</table>

Preprocessing

$H' = (H - \mu)/\sqrt{(\sigma^2 + c)}, \text{where}$

$H$ is the historical horizon

$\mu$ is the mean

$\sigma$ is the variance

$c$ ensures numerical stability

StandardScaler

$z = (x - \mu)/\sigma, \text{where}$

$x$ is the sample

$\mu$ is the mean

$\sigma$ is the standard deviation

**What are they actually doing?**

<div style="columns:2">

### Preprocessing

$$H' = (H - \mu)/\sqrt{(\sigma^2 + c)}, \text{where}$$

$H$ is the historical horizon

$\mu$ is the mean

$\sigma$ is the variance

$c$ ensures numerical stability

- We know that $\sqrt{\sigma^2} = \sigma$

### StandardScaler

$$z = (x - \mu)/\sigma, \text{where}$$

$x$ is the sample

$\mu$ is the mean

$\sigma$ is the standard deviation

</div>

**What are they actually doing?**

Preprocessing

$$H' = (H - \mu)/\sqrt{(\sigma^2 + c)}, \text{where}$$

$H$ is the historical horizon

$\mu$ is the mean

$\sigma$ is the variance

$c$ ensures numerical stability

StandardScaler

$$z = (x - \mu)/\sigma, \text{where}$$

$x$ is the sample

$\mu$ is the mean

$\sigma$ is the standard deviation

- We know that $\sqrt{\sigma^2} = \sigma$
- Essentially same formula, except constant

**What are they actually doing?**

<div style="display: flex;">
<div>

### Preprocessing

$H' = (H - \mu)/\sqrt{(\sigma^2 + c)},$ where

$H$ is the historical horizon

$\mu$ is the mean

$\sigma$ is the variance

$c$ ensures numerical stability

</div>
<div>

### StandardScaler

$z = (x - \mu)/\sigma,$ where

$x$ is the sample

$\mu$ is the mean

$\sigma$ is the standard deviation

</div>
</div>

- We know that $\sqrt{\sigma^2} = \sigma$
- Essentially same formula, except constant
- Fit on training data, normalize entire dataset $\rightarrow$ global normalization

**What are they actually doing?**

<table>
<tr><td align="center">Preprocessing</td><td align="center">StandardScaler</td></tr>
</table>

$H' = (H - \mu)/\sqrt{(\sigma^2 + c)}, \text{where}$

$H$ is the historical horizon

$\mu$ is the mean

$\sigma$ is the variance

$c$ ensures numerical stability

$z = (x - \mu)/\sigma, \text{where}$

$x$ is the sample

$\mu$ is the mean

$\sigma$ is the standard deviation

- We know that $\sqrt{\sigma^2} = \sigma$
- Essentially same formula, except constant
- Fit on training data, normalize entire dataset $\rightarrow$ global normalization
- None of the stated benefits of instance normalization
  - ▸ Mitigate internal covariate shift
  - ▸ Grasp intricate temporal dynamics in TS

## Architecture

Upper part $\rightarrow$ dynamic correlation

Lower part $\rightarrow$ normalized data

Output $\rightarrow$ enriched input features

**Normalized Data**

- Normalized as described earlier
  - ‣ Not what the paper actually states

**Normalized Data**

- Normalized as described earlier
  - ‣ Not what the paper actually states

**Patches**

- $H'$ is split into $p$ patches
- Stride $S$
- Size $T$
- If $S \geq T$ patches are disjoint
- If $S < T$ patches overlap
  - ‣ Common elements for adjacent patches

**AME**

- Provides local memory embedding
  - ‣ These are learnable parameters
- Consistant local memory for patch $P_i$
- Matrix product of $E_i \otimes E_i^T$
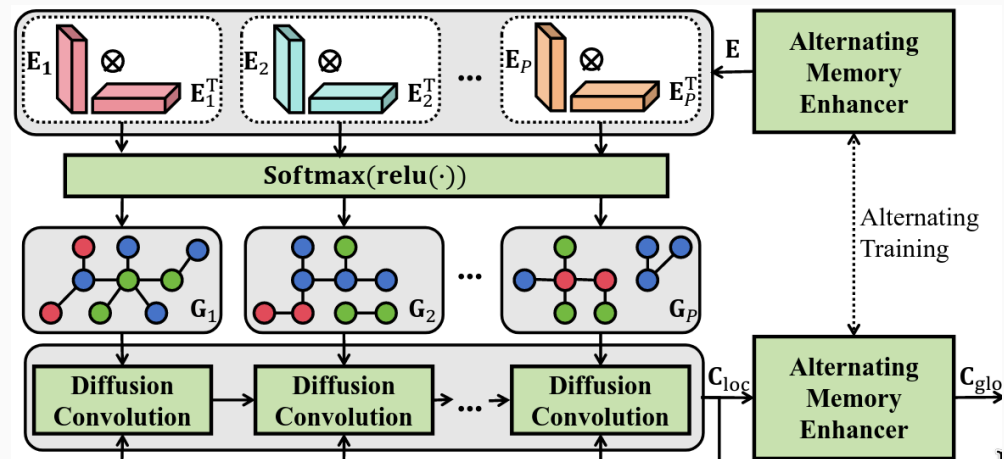  - ‣ Similarity matrix for variables in $P_i$

## AME

- Provides local memory embedding
  - ▸ These are learnable parameters
- Consistant local memory for patch $P_i$
- Matrix product of $E_i \otimes E_i^T$
  - ▸ Similarity matrix for variables in $P_i$

## ReLU + Softmax

- ReLU eliminates negative values
  - ▸ Removes negative correlations
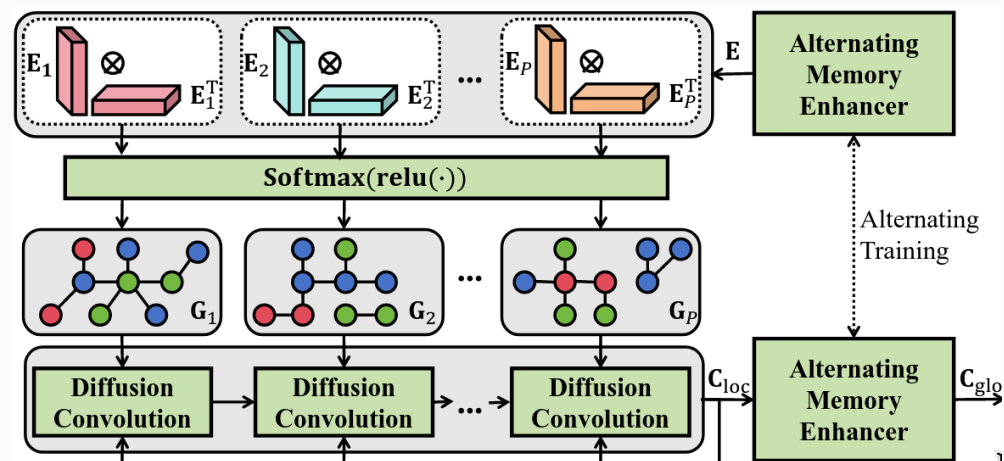- Softmax scales into influence scores

## AME

- Provides local memory embedding
  - ▸ These are learnable parameters
- Consistant local memory for patch $P_i$
- Matrix product of $E_i \otimes E_i^T$
  - ▸ Similarity matrix for variables in $P_i$

## ReLU + Softmax

- ReLU eliminates negative values
  - ▸ Removes negative correlations
- Softmax scales into influence scores

## Graph

- Translates influence scores into graph
- Captures connection between variables
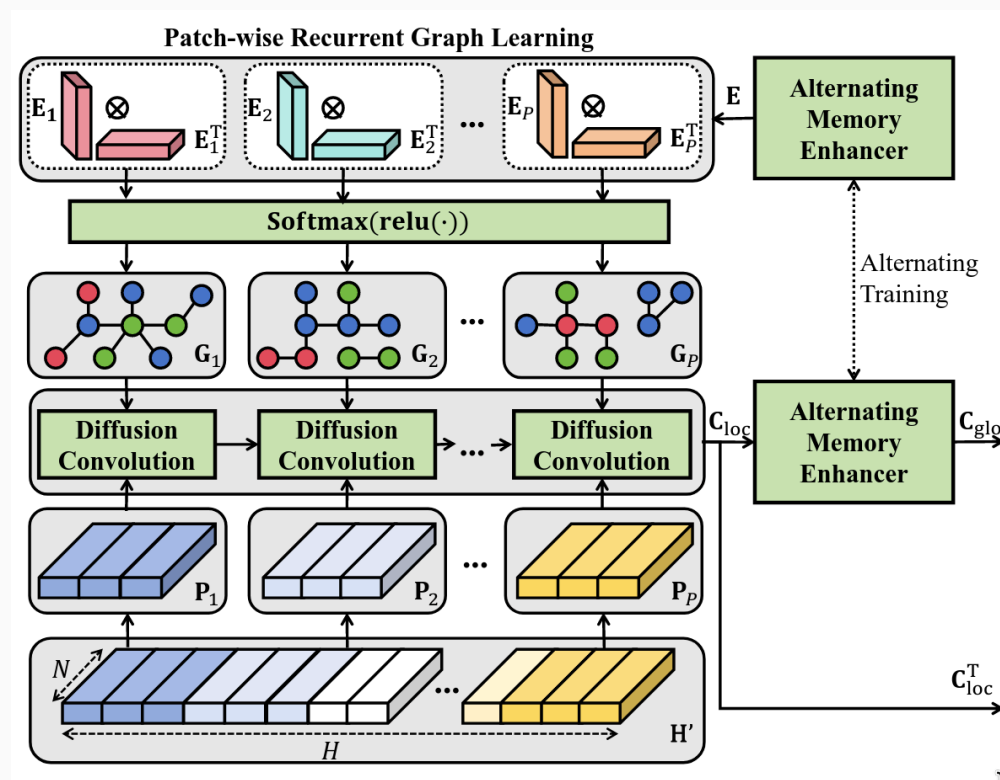  - ▸ Dynamic correlations

**Diffusion Convolution**

- Normalized data is adjusted based on connections in graph
- Numeric values "diffuse" into neighbours
  ‣ Not only immediate neighbours
- Spatially relates data based on connections

**Gated Recurrent Unit**

- Forwards information from $P_i$ to $P_{i+1}$
- Temporally relates data in a sequence

**Output**

- Input features enriched with local information
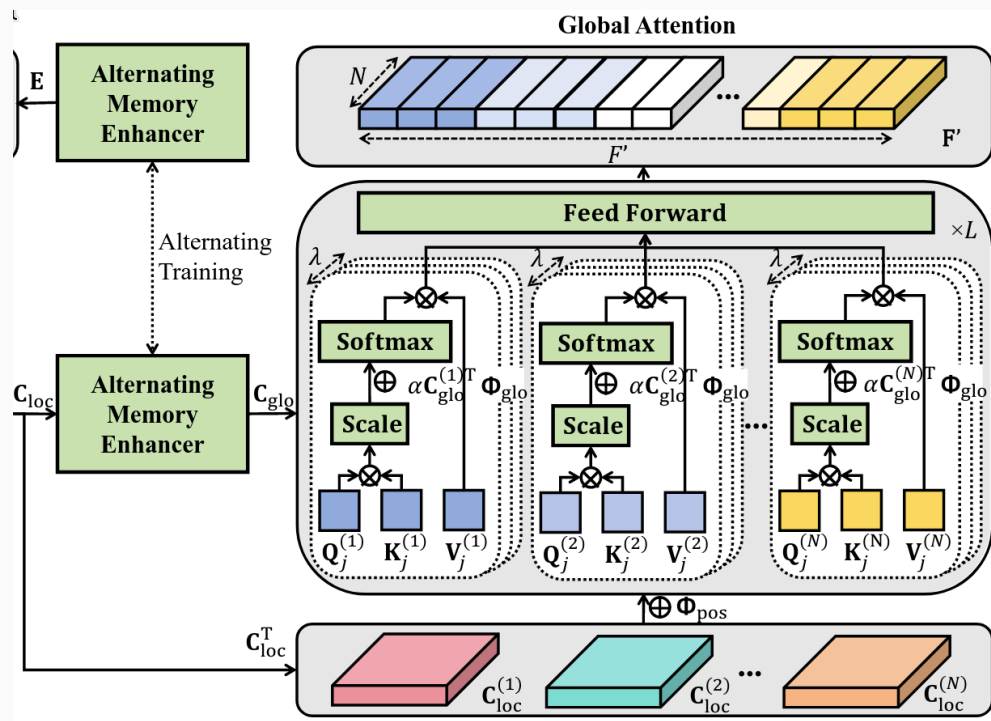- Spatial $\rightarrow$ dynamic correlations
- Temporal $\rightarrow$ GRU

**Motivation**

- Patch-wise correlations are sensitive
  - ▸ Outliers dominate
- Constrain locally enriched features
  - ▸ Mitigate disrupted correlations

**Input**

- Transpose locally enriched features
  - ▸ Isolate variables
  - ▸ Diffusion earlier
- Linear transformation
  - ▸ Positional encoding
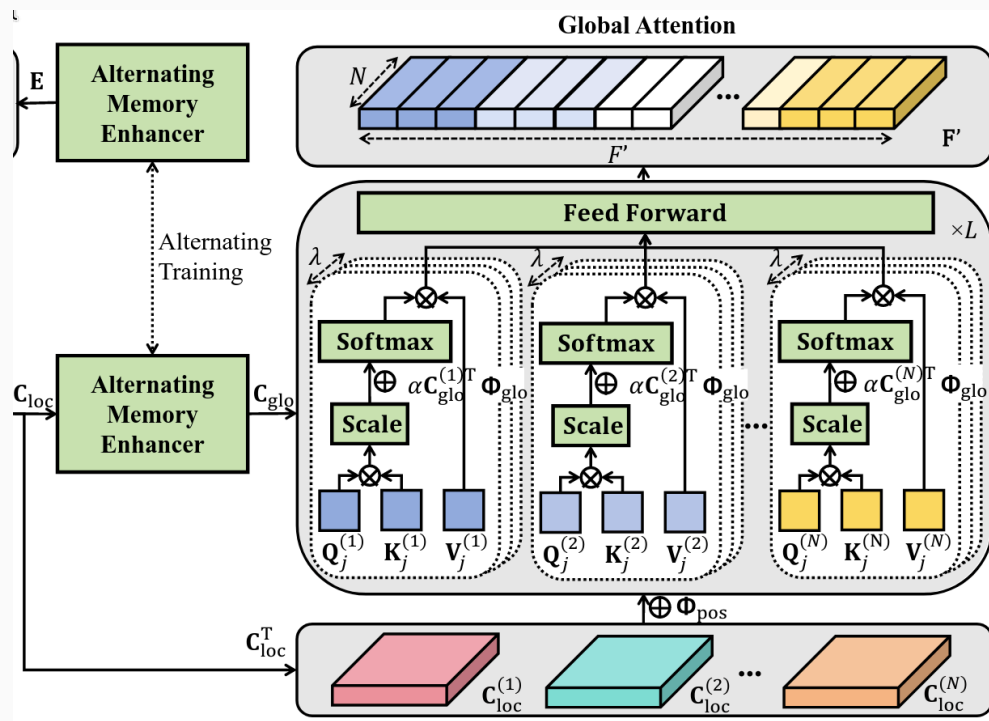- Converted to Q, K, V matrices
  - ▸ Learnable parameters

**Attention**

- Relatively conventional implementation
  - ▸ Query and Key to find importance
  - ▸ Weight Value by importance
- Global information is new
- Adding global information after softmax
  - ▸ Bias probabilities
  - ▸ Global information affects parameters

**Output**

- The final "representation" of data
- **F'** is not a forecast
  - ▸ Final feature representation
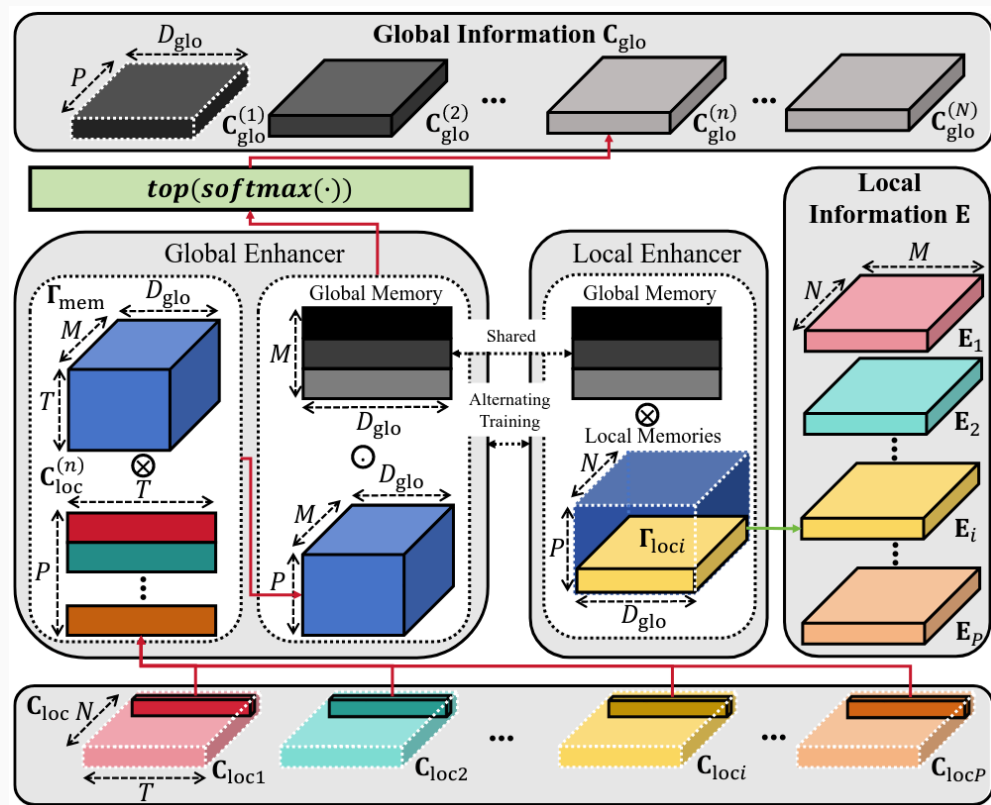- Linear layer maps to forecasting horizon

**Overview**

- Input
  - ‣ Locally correlated features
- Outputs
  - ‣ Local information $E$
  - ‣ Global information $C_{\text{glo}}$
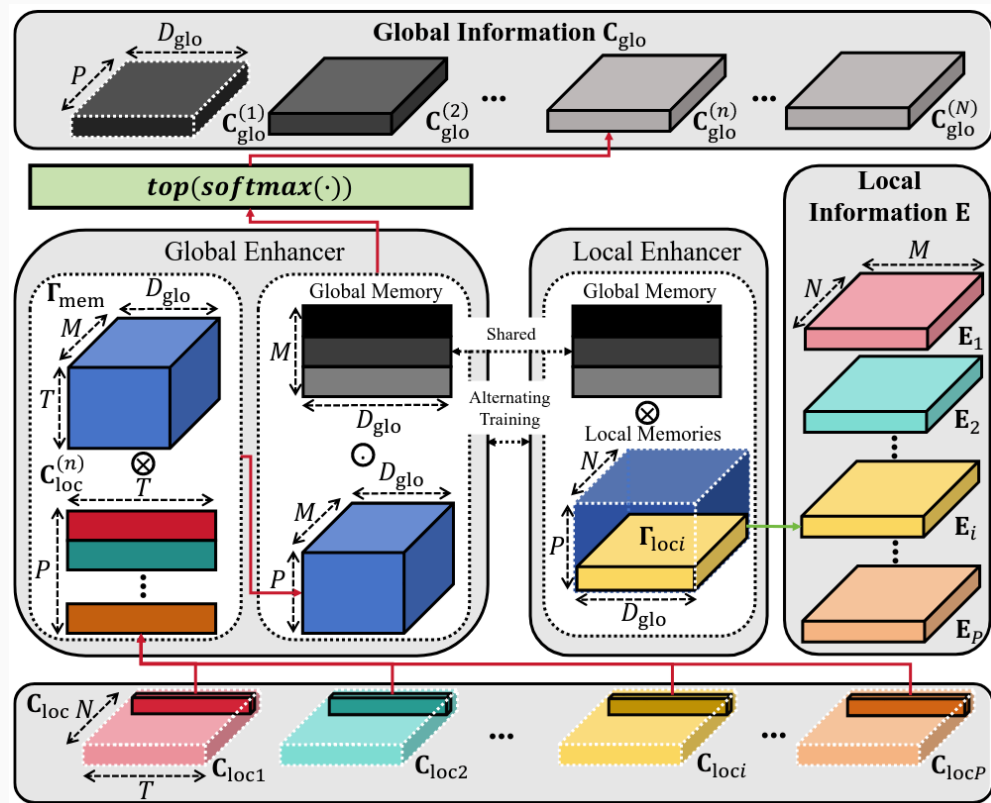- Shared global memory

**Hyperparameters**

- $M \rightarrow$ number of high level patterns
  - ‣ Spikes, seasons, stable
- $D_{\text{glo}} \rightarrow$ richness of patterns

## Local Enhancer

- Local memory regions $\Gamma_{\text{loc}i}$
  - ▸ One for each patch
- $P_i \longleftrightarrow \Gamma_{\text{loc}i} \rightarrow \Gamma_{\text{loc}i} \longleftrightarrow E_i$
- Not directly identical
  - ▸ $E$ contains global memories
  - ▸ Defined by $C_{\text{loc}}$
- Memories are **not** information

## Global Enhancer

- Learns from locally correlated features
- $\Gamma_{\mathrm{mem}}$ is a large trainable tensor
  - ‣ Produces inquiry tensor
  - ‣ Recognizes patterns in data
    - – The $M$ high level patterns
- Inquiry tensor
  - ‣ Prevalence of patterns in local data
  - ‣ Similarity scores with global memory
- Probability distribution
  - ‣ Importance of pattern
- Top $k$ most important patterns
  - ‣ Stored in $C_{\mathrm{glo}}$
  - ‣ Scaled based on importance
    - – Weighted sum

**Alternating Training**

- Local information $E$ requires
  - ‣ Local memories
  - ‣ Global memories
- Updating both simultaneously
  - ‣ Unstable training
  - ‣ Issues converging
- LE and GE alternate training
  - ‣ Split adjustment of memories

**LE Training**

- Local memories > global memories
  - ‣ More parameters $\rightarrow$ longer convergence
- Balance convergence
  - ‣ Different learning rates
  - ‣ LE training more

## Alternating Training

- Local information $E$ requires
  - ▸ Local memories
  - ▸ Global memories
- Updating both simultaneously
  - ▸ Unstable training
  - ▸ Issues converging
- LE and GE alternate training
  - ▸ Split adjustment of memories

## LE Training

- Local memories > global memories
  - ▸ More parameters $\rightarrow$ longer convergence
- Balance convergence
  - ▸ Different learning rates
  - ▸ LE training more

---

**Algorithm 1** AME alternating training

---

**Input:** Historical horizon and ground truth $\mathbf{H}, \mathbf{F}$; local and global memories $\Gamma_{\text{loc}}, \Gamma_{\text{glo}}$; local training step $\epsilon$; learning rates $\eta_{\text{loc}}, \eta_{\text{glo}}$ for local and global enhancers

**Output:** Local and global information $\mathbf{E}, \mathbf{C}_{\text{glo}}$; learned local and global memories $\Gamma_{\text{loc}}$, $\Gamma_{\text{glo}}$, tensor $\Gamma_{\text{mem}}$, and bias $\mathbf{b}_{\text{mem}}$

1: *Initialisation*: Initializing local and global memories $\Gamma_{\text{loc}}, \Gamma_{\text{glo}}$, tensor $\Gamma_{\text{mem}}$, and bias $\mathbf{b}_{\text{mem}}$ randomly

2: **while** $\Gamma_{\text{loc}}, \Gamma_{\text{glo}}, \Gamma_{\text{mem}}$, and $\mathbf{b}_{\text{mem}}$ are not converged **do**

3:     **for** iteration = 0 to $\epsilon$ **do**

4:         $\mathbf{H}' \leftarrow \text{Preprocessing}(\mathbf{H})$

5:         $\mathbf{E} \leftarrow \mathcal{A}_{\text{loc}}(\Gamma_{\text{loc}}, \Gamma_{\text{glo}})$

6:         $\mathbf{C}_{\text{loc}} \leftarrow \mathcal{G}_{\Theta}(\mathbf{H}', \mathbf{E})$

7:         $\mathbf{C}_{\text{glo}} \leftarrow \mathcal{A}_{\text{glo}}(\mathbf{C}_{\text{loc}}, \Gamma_{\text{glo}})$

8:         $\mathbf{F}' \leftarrow \mathcal{T}_{\Phi}(\mathbf{C}_{\text{loc}}, \mathbf{C}_{\text{glo}})$

9:         $\hat{\mathbf{F}} \leftarrow \text{LinearHead}(\mathbf{F}')$

10:         $\Gamma_{\text{loc}} \leftarrow \Gamma_{\text{loc}} - \eta_{loc} \nabla_{\Gamma_{\text{loc}}} \mathcal{L}(\hat{\mathbf{F}}, \mathbf{F})$

11:     **end for**

12:     $\mathbf{H}' \leftarrow \text{Preprocessing}(\mathbf{H})$

13:     $\mathbf{E} \leftarrow \mathcal{A}_{\text{loc}}(\Gamma_{\text{loc}}, \Gamma_{\text{glo}})$

14:     $\mathbf{C}_{\text{loc}} \leftarrow \mathcal{G}_{\Theta}(\mathbf{H}', \mathbf{E})$

15:     $\mathbf{C}_{\text{glo}} \leftarrow \mathcal{A}_{\text{glo}}(\mathbf{C}_{\text{loc}}, \Gamma_{\text{glo}})$

16:     $\mathbf{F}' \leftarrow \mathcal{T}_{\Phi}(\mathbf{C}_{\text{loc}}, \mathbf{C}_{\text{glo}})$

17:     $\hat{\mathbf{F}} \leftarrow \text{LinearHead}(\mathbf{F}')$

18:     $\Gamma_{\text{glo}} \leftarrow \Gamma_{\text{glo}} - \eta_{glo} \nabla_{\Gamma_{\text{glo}}} \mathcal{L}(\hat{\mathbf{F}}, \mathbf{F})$

19:     $\Gamma_{\text{mem}} \leftarrow \Gamma_{\text{mem}} - \eta_{glo} \nabla_{\Gamma_{\text{mem}}} \mathcal{L}(\hat{\mathbf{F}}, \mathbf{F})$,

20:     $\mathbf{b}_{\text{mem}} \leftarrow \mathbf{b}_{\text{mem}} - \eta_{glo} \nabla_{\mathbf{b}_{\text{mem}}} \mathcal{L}(\hat{\mathbf{F}}, \mathbf{F})$

21: **end while**

# 3. Results

# 4. Critique

**Preprocessing**

- As mentioned earlier
- Unconventional notation
- Obscures details

**Inconsistencies**

- $C_{\text{glo}}$ is global memory
- $C_{\text{loc}}$ is locally correlated features
- $E$ is local memory

## Preprocessing

- As mentioned earlier
- Unconventional notation
- Obscures details

## Inconsistencies

- $C_{\text{glo}}$ is global memory
- $C_{\text{loc}}$ is locally correlated features
- $E$ is local memory

## Symbol Reuse

- $\mathbf{F}$ is the ground truth
- $F$ is the dimensionality of $\mathbf{F}$
- $\mathbf{F'}$ is the encoding output
- $F'$ is the dimensionality of $\mathbf{F'}$
- Confusing statements and diagrams

**Preprocessing**

- As mentioned earlier
- Unconventional notation
- Obscures details

**Inconsistencies**

- $C_{\mathrm{glo}}$ is global memory
- $C_{\mathrm{loc}}$ is locally correlated features
- $E$ is local memory

**Symbol Reuse**

- $\mathbf{F}$ is the ground truth
- $F$ is the dimensionality of $\mathbf{F}$
- $\mathbf{F'}$ is the encoding output
- $F'$ is the dimensionality of $\mathbf{F'}$
- Confusing statements and diagrams

$\mathbf{F'} \in \mathbb{R}^{F' \times N}$, where $F'$ is the temporal dimension of the representation.

## Preprocessing

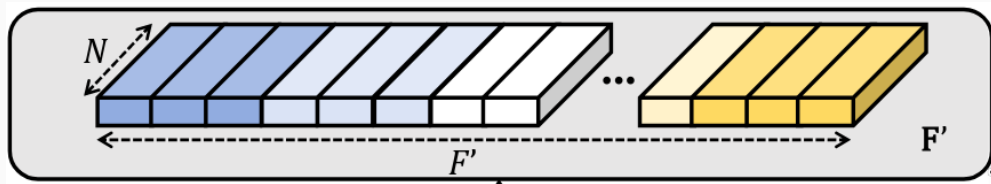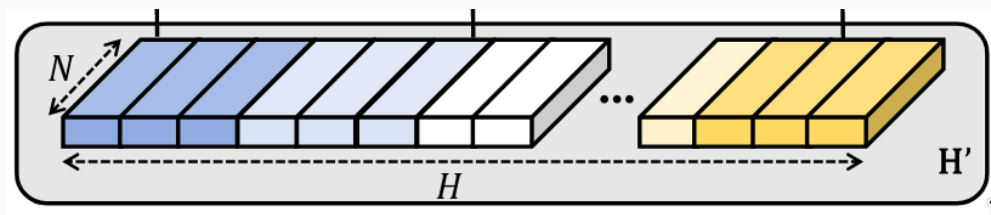- As mentioned earlier
- Unconventional notation
- Obscures details

## Inconsistencies

- $C_{\text{glo}}$ is global memory
- $C_{\text{loc}}$ is locally correlated features
- $E$ is local memory

## Symbol Reuse

- $\mathbf{F}$ is the ground truth
- $F$ is the dimensionality of $\mathbf{F}$
- $\mathbf{F}'$ is the encoding output
- $F'$ is the dimensionality of $\mathbf{F}'$
- Confusing statements and diagrams

$\mathbf{F}' \in \mathbb{R}^{F' \times N}$, where $F'$ is the temporal dimension of the representation.

# 5. Praise

## Colors

- Help understanding and data flow
  - ‣ Preprocessing → final encoding
  - ‣ Minor inconsistencies
    - – Attention

## Dimensionality

- Squares → 2-dimensional
- Cubes → 3-dimensional
- Transposed → lying down
- Slices of shapes
  - ‣ $M$ slices of global memory
  - ‣ $P$ slices of local memory