

Memformer

A Memory Guided Transformer for Time Series Forecasting

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

February 2025

Proceedings of the VLDB Endowment, Volume 18, Issue 2

Presented by **Andreas Gottschalk Krath**

1. Introduction



1.1 Motivation

Forecasting

- Predicting the future
 - Allows preparation

1.1 Motivation

Forecasting

- Predicting the future
 - Allows preparation
- Long term forecasting?
 - Obviously more difficult than short term
 - Time constrained tasks

1.1 Motivation

Long Term Forecasting

- What defines long term?

1.1 Motivation

Long Term Forecasting

- What defines long term?
- Historical horizon
- Forecasting horizon

1.1 Motivation

Long Term Forecasting

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

Variable Correlation

- Complex systems have many variables
 - These relate to each other
- These impact forecasting accuracy
 - Patterns in the data

1.1 Motivation

Dynamic Correlations

- Are variable correlations stable over time?
 - No

1.1 Motivation

Dynamic Correlations

- Are variable correlations stable over time?
 - No
- Correlations are dynamic over time
 - Seasons
 - Sensor drift

1.1 Motivation

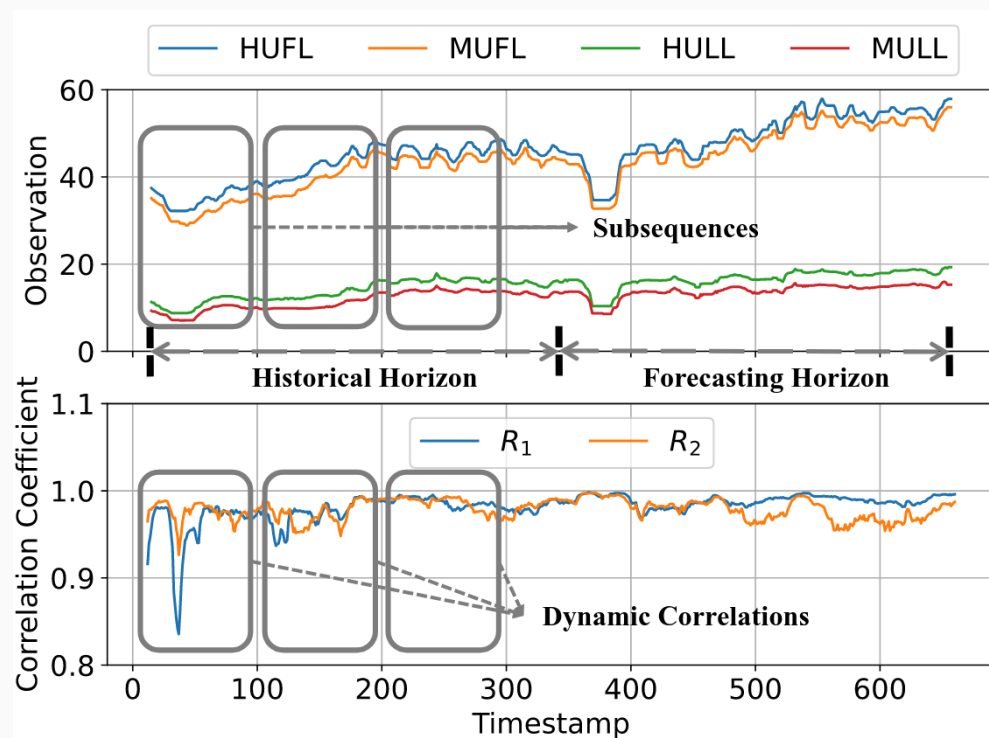
Dynamic Correlations

- 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

1.1 Motivation

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 $R_2 = 0.990$.

1.1 Motivation

Disrupted Correlations

- System errors
- External influence

1.1 Motivation

Disrupted Correlations

- System errors
- External influence
- What happens with outliers?
 - Affect correlation \rightarrow accuracy

1.1 Motivation

Disrupted Correlations

- System errors
- External influence
- What happens with outliers?
 - Affect correlation \rightarrow accuracy
- Many models are sensitive to outliers
 - Numeric difference dominates training

1.1 Motivation

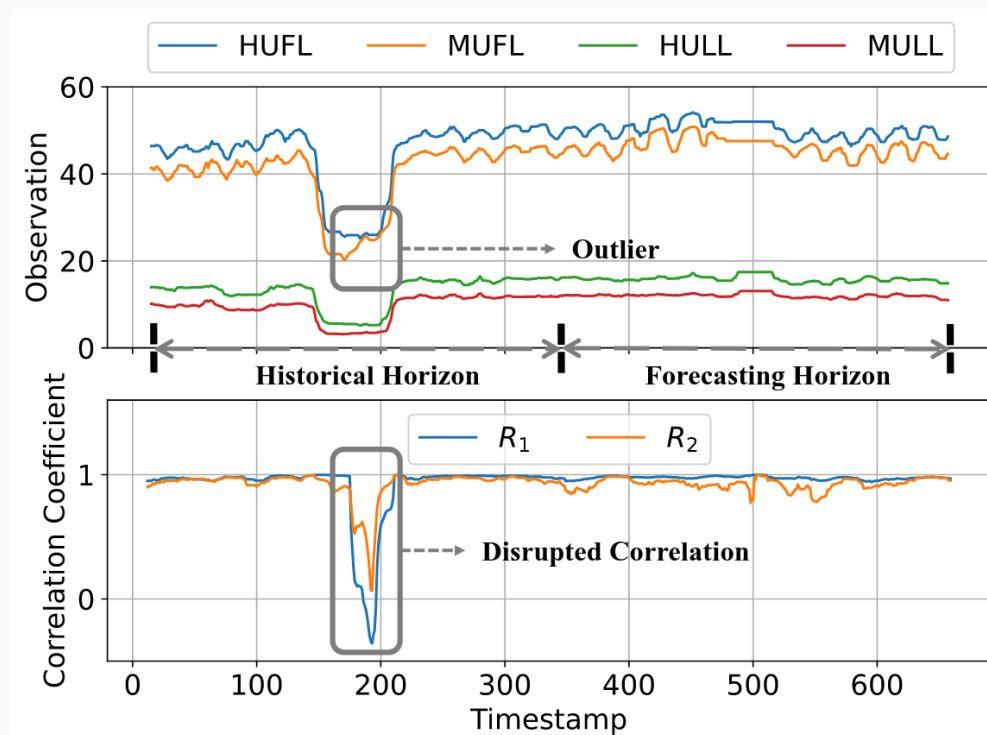
Disrupted Correlations

- 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

1.1 Motivation

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 $R_2 = 0.963$.

1.2 Problem

Challenge 1

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

1.2 Problem

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

1.3 Contributions

Memformer

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

1.3 Contributions

Memformer

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

Alternating Memory Enhancer

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

1.3 Contributions

Memformer

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

Alternating Memory Enhancer

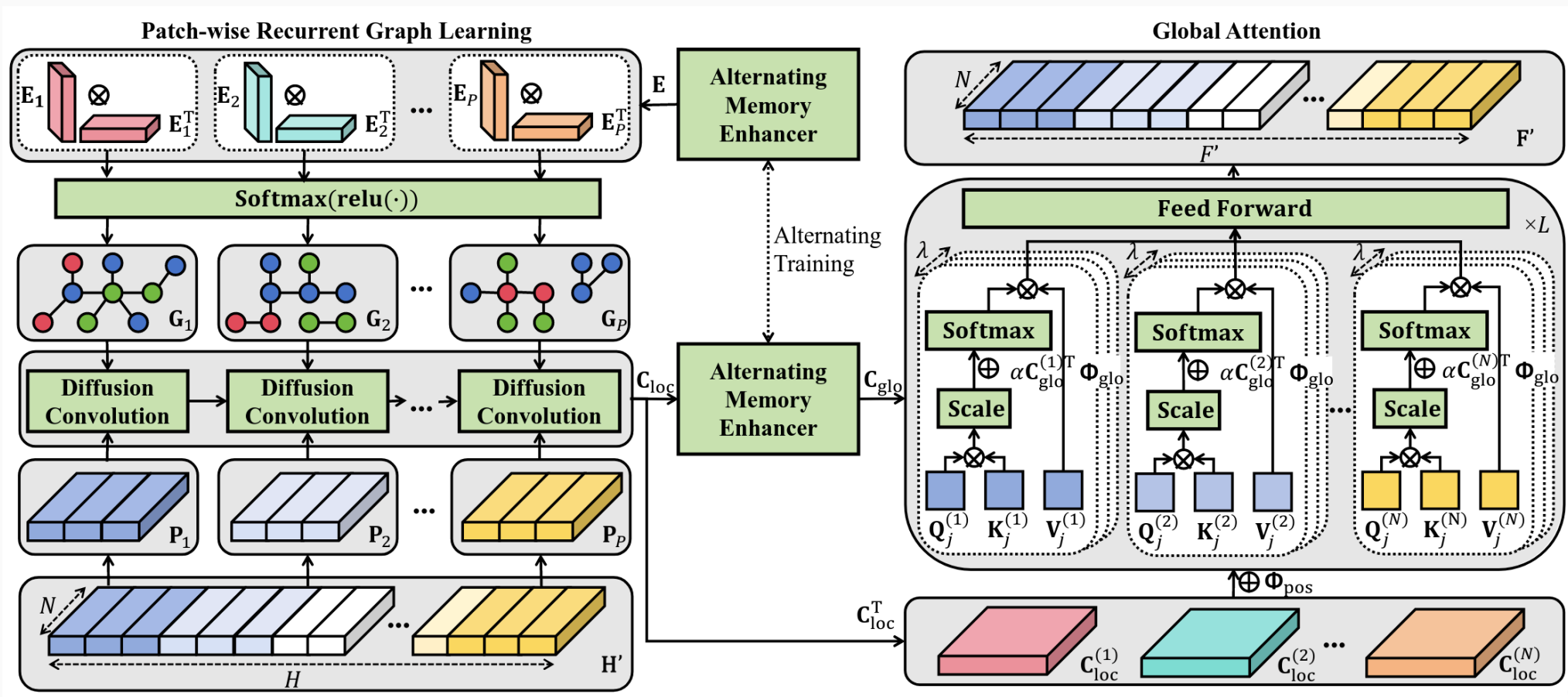
- Memory network
- Associates local and global information
- Addresses challenge 2

Experiments

- Proof

2. Methodology

2.1 Overview



1

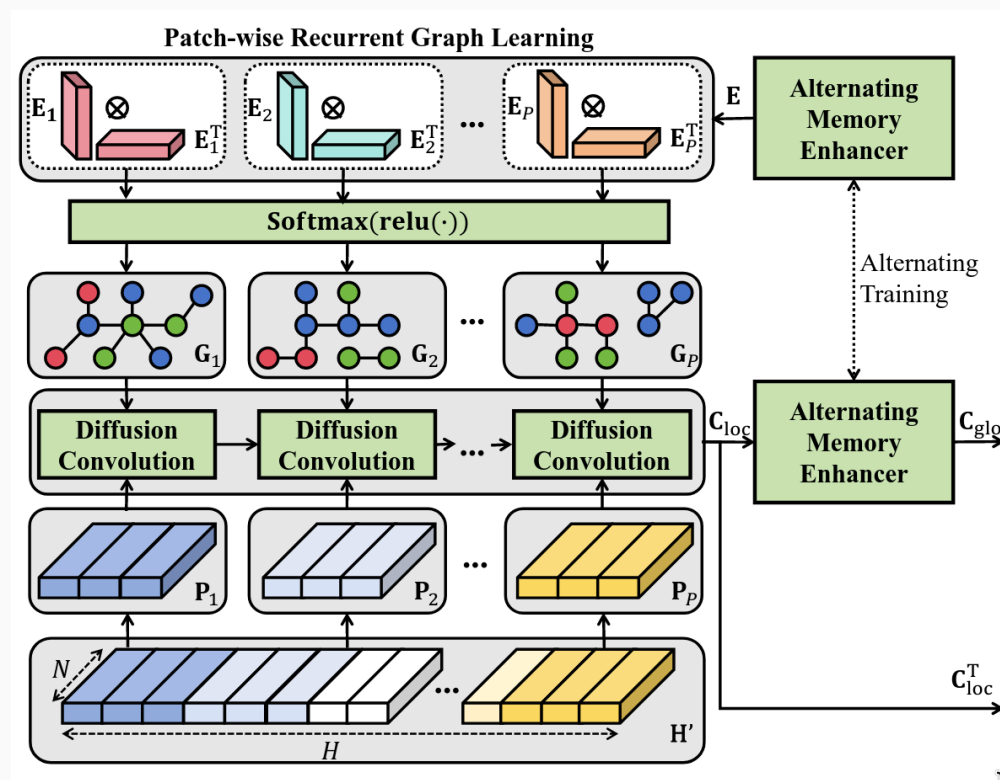
2.2 Patch-wise Recurrent Graph Learning

Architecture

Upper part \rightarrow dynamic correlation

Lower part \rightarrow normalized data

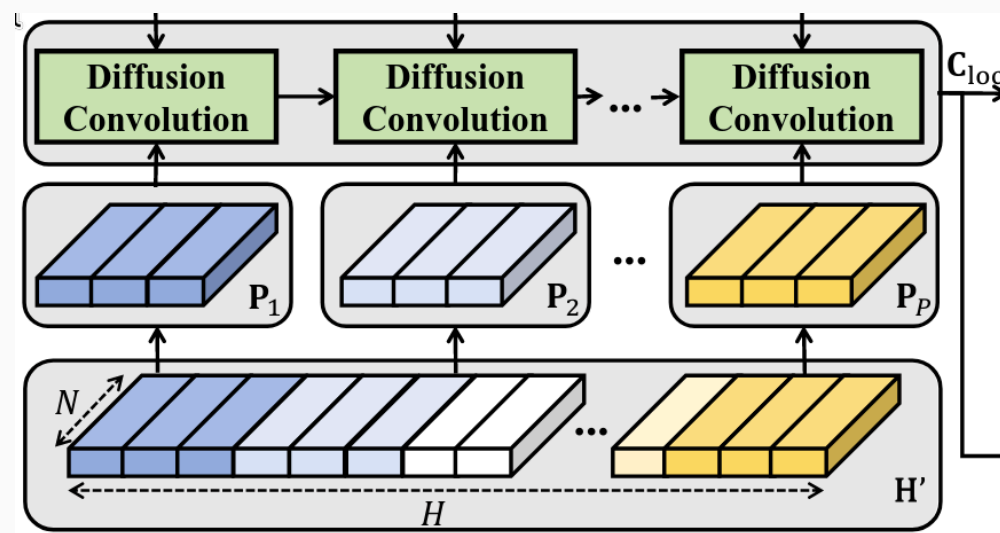
Output \rightarrow enriched input features



2.2 Patch-wise Recurrent Graph Learning

Normalized Data

- Instance normalization



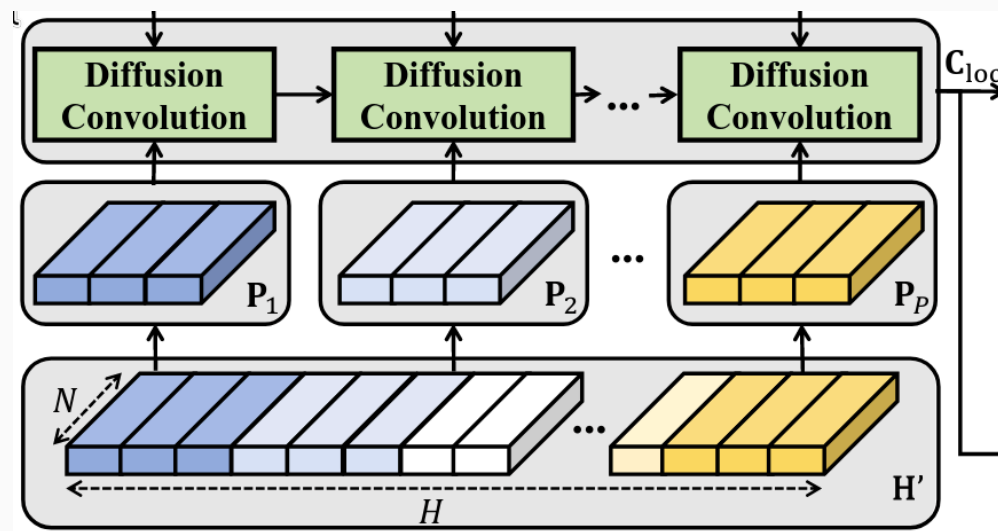
2.2 Patch-wise Recurrent Graph Learning

Normalized Data

- Instance normalization

Patches

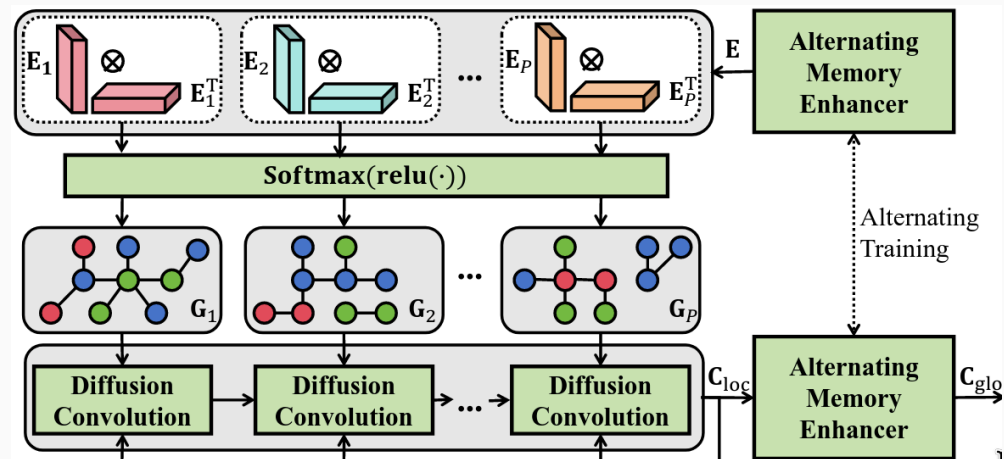
- H' is split into p patches
- Group temporally related data



2.2 Patch-wise Recurrent Graph Learning

AME

- Provides local information
 - These are learnable parameters
- Consistent local information for patch P_i
- Matrix product of $E_i \otimes E_i^T$
 - Similarity matrix for variables in P_i



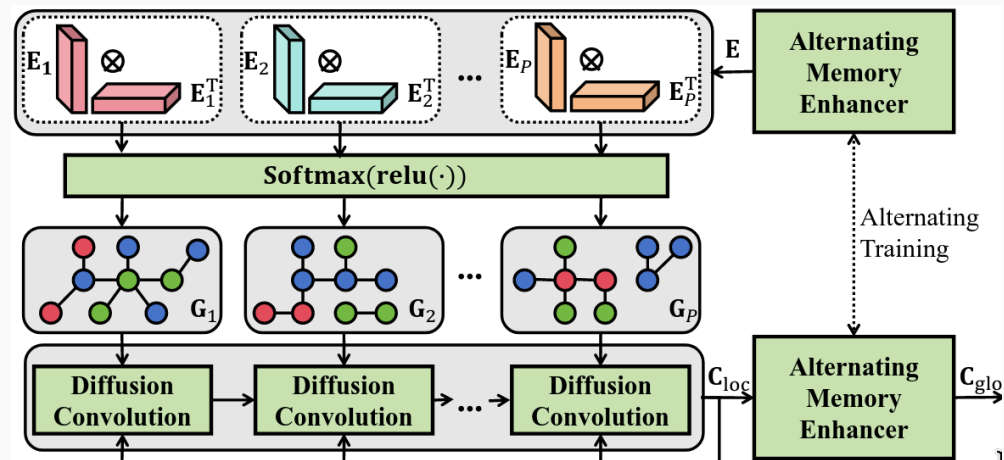
2.2 Patch-wise Recurrent Graph Learning

AME

- Provides local information
 - These are learnable parameters
- Consistent local information 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



2.2 Patch-wise Recurrent Graph Learning

AME

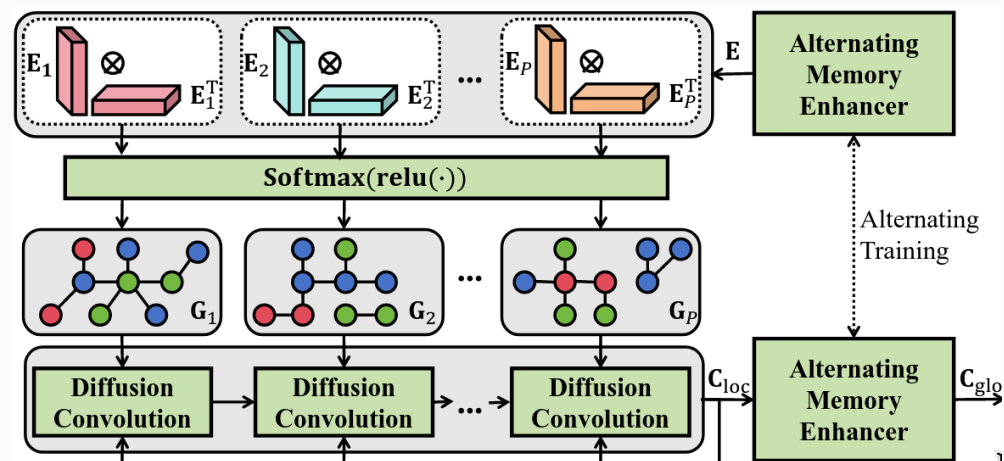
- Provides local information
 - These are learnable parameters
- Consistent local information for patch P_i
- Matrix product of $E_i \otimes E_i^T$
 - Similarity matrix for variables in P_i

ReLU + Softmax

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

Graph

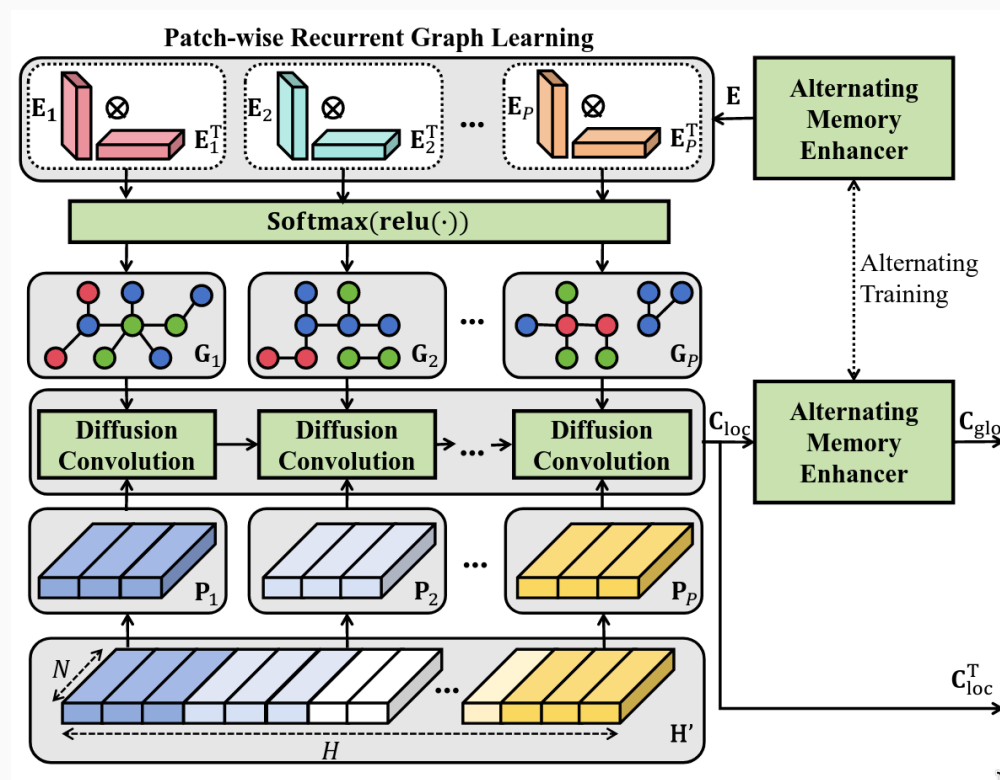
- Translates influence scores into graph
- Captures connection between variables
 - Dynamic correlations



2.2 Patch-wise Recurrent Graph Learning

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



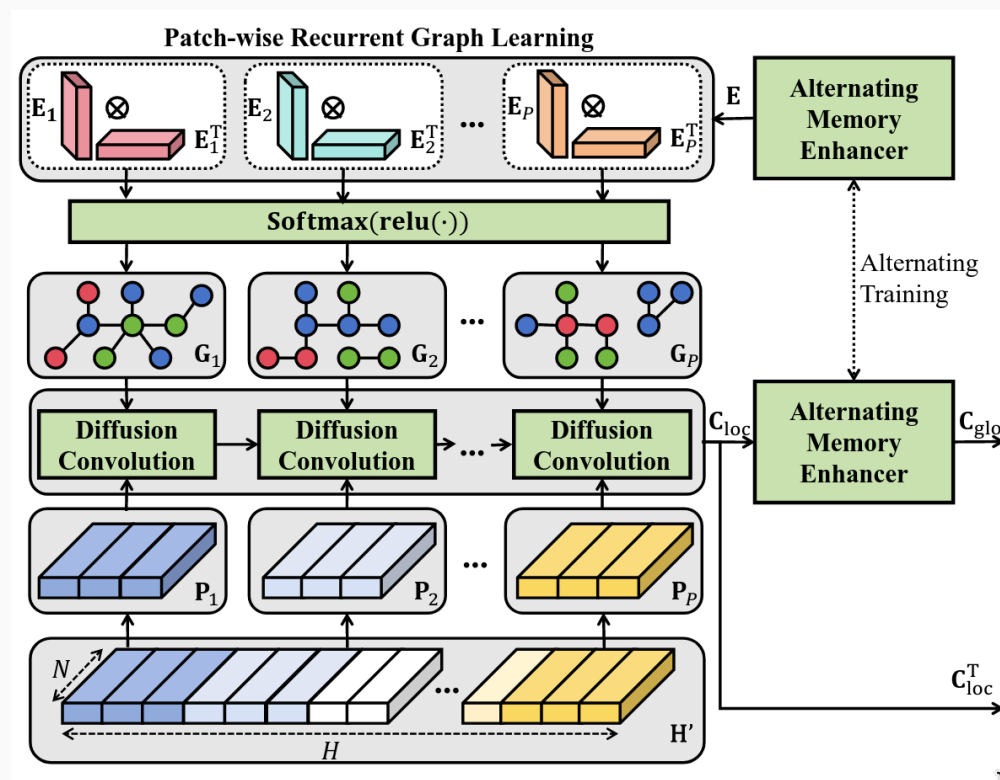
2.2 Patch-wise Recurrent Graph Learning

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



2.2 Patch-wise Recurrent Graph Learning

Diffusion Convolution

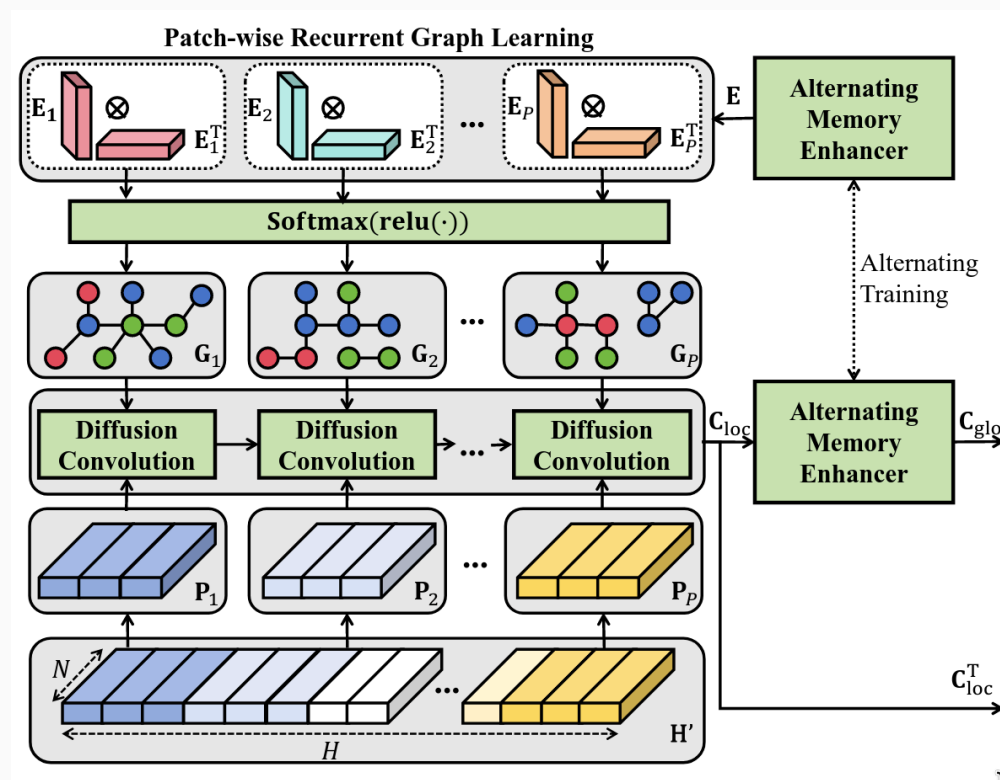
- Normalized data is adjusted based on connections in graph
- Numeric values “diffuse” into neighbours
 - Not only immediate neighbours
- Spatially relates data based on connections

Gated Recurrent Unit

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

Output

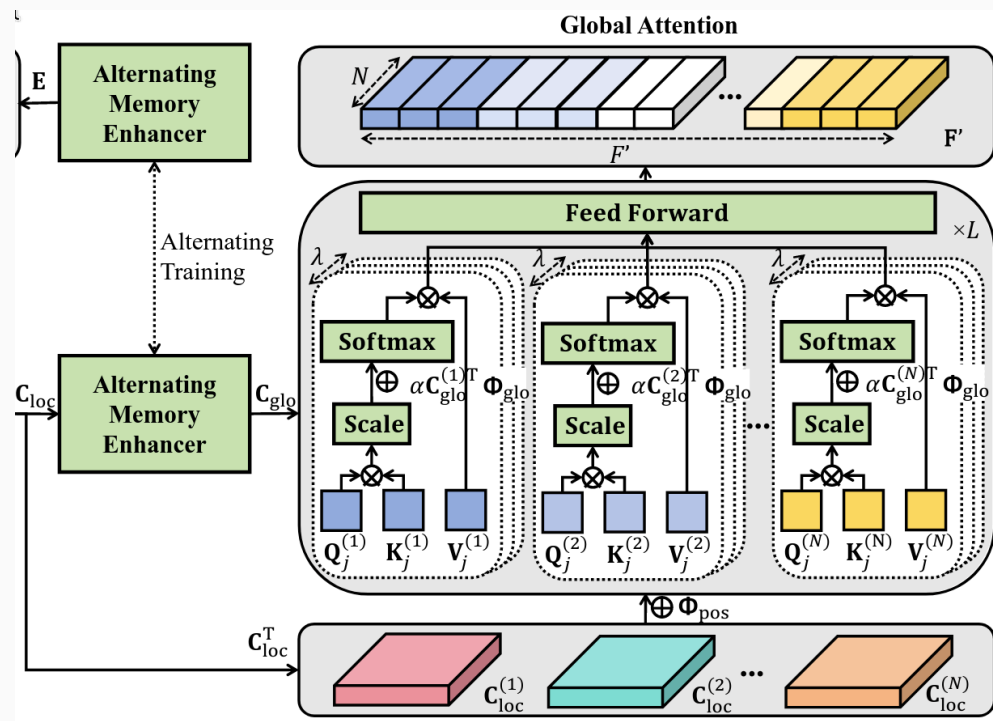
- Input features enriched with local information
- Spatial \rightarrow dynamic correlations
- Temporal \rightarrow GRU



2.3 Global Attention

Motivation

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



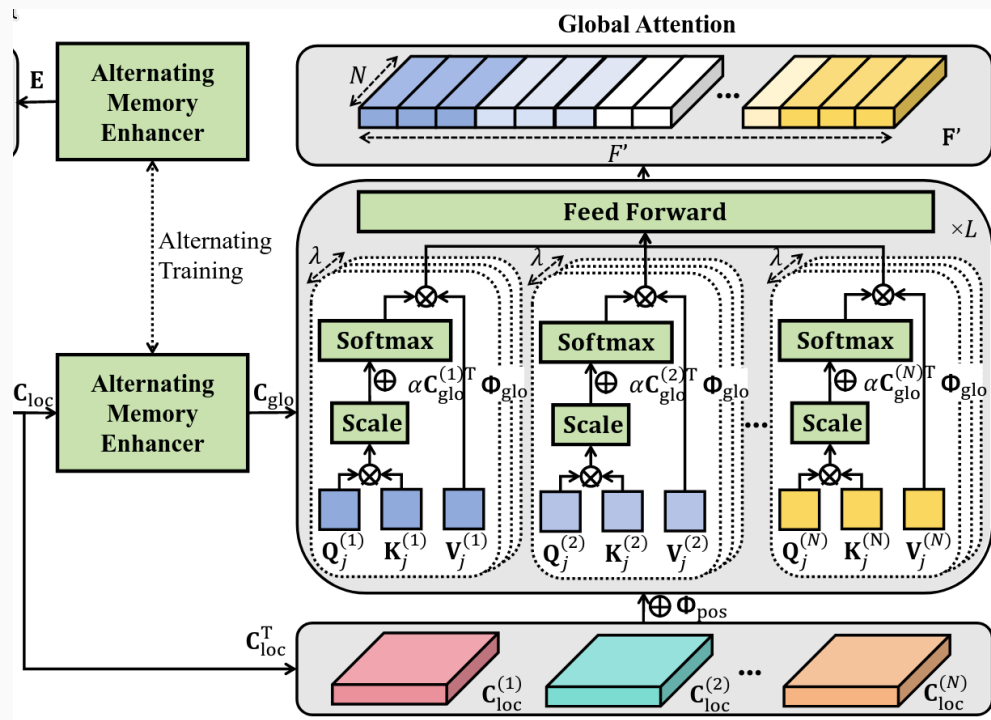
2.3 Global Attention

Motivation

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

Input

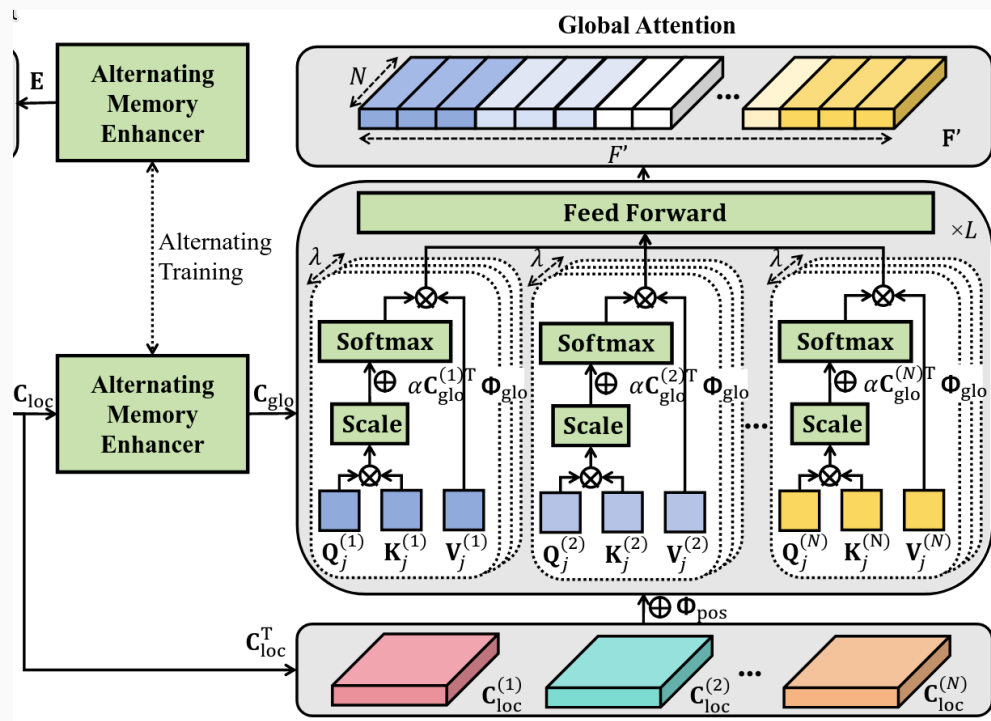
- Transpose locally enriched features
 - Isolate variables
 - Diffusion earlier
- Converted to Q, K, V matrices
 - Learnable parameters



2.3 Global Attention

Attention

- Relatively conventional implementation
- Global information is new
- Adding global information after softmax
 - Bias probabilities



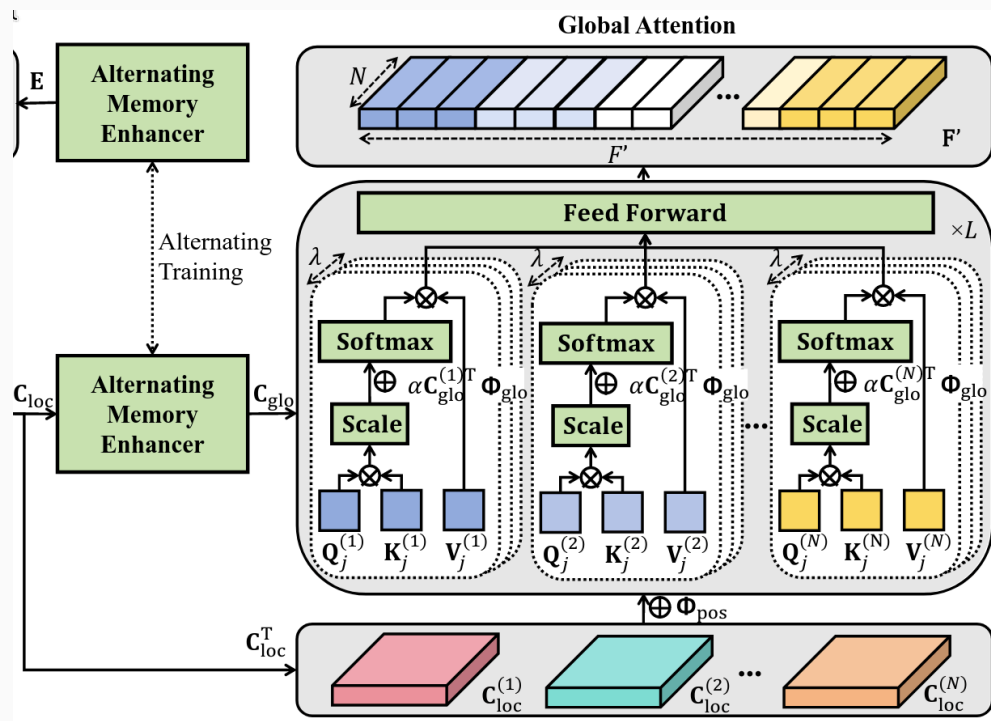
2.3 Global Attention

Attention

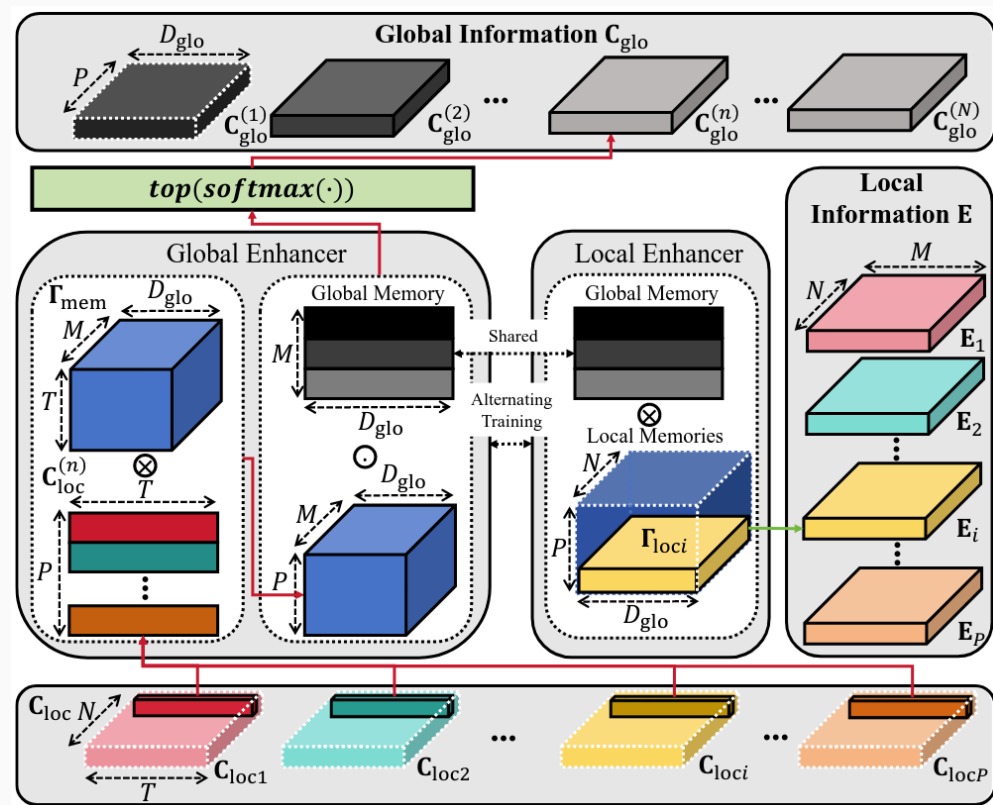
- Relatively conventional implementation
- Global information is new
- Adding global information after softmax
 - Bias probabilities

Output

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



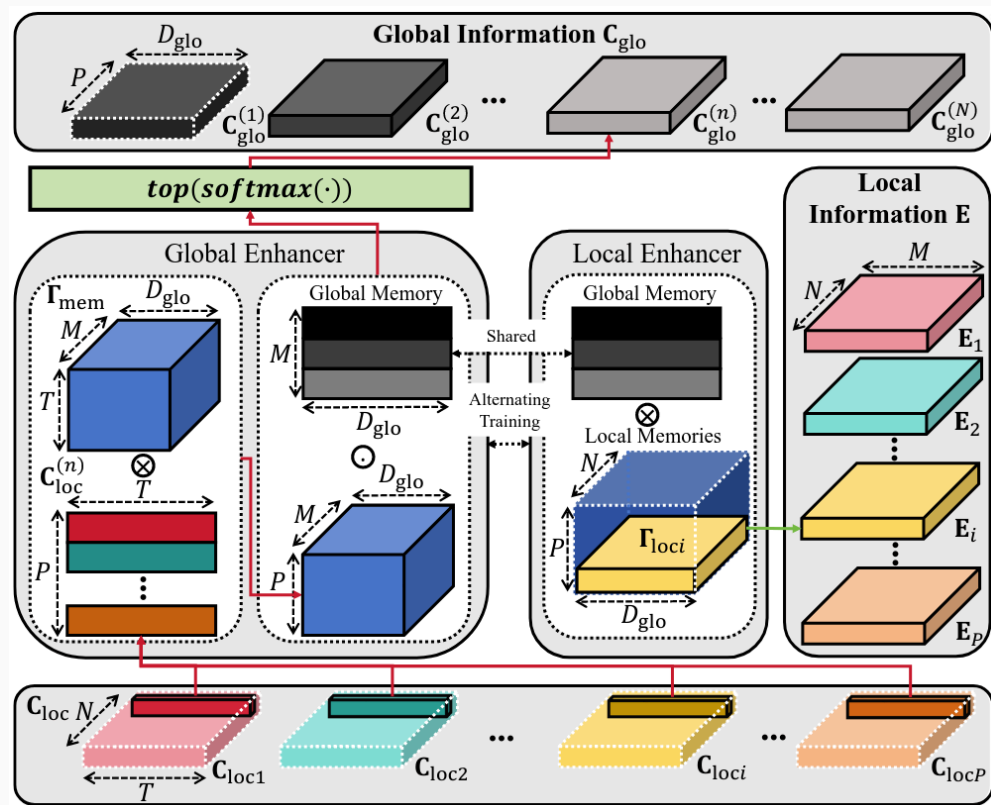
2.4 Alternating Memory Enhancer



2.4 Alternating Memory Enhancer

Overview

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



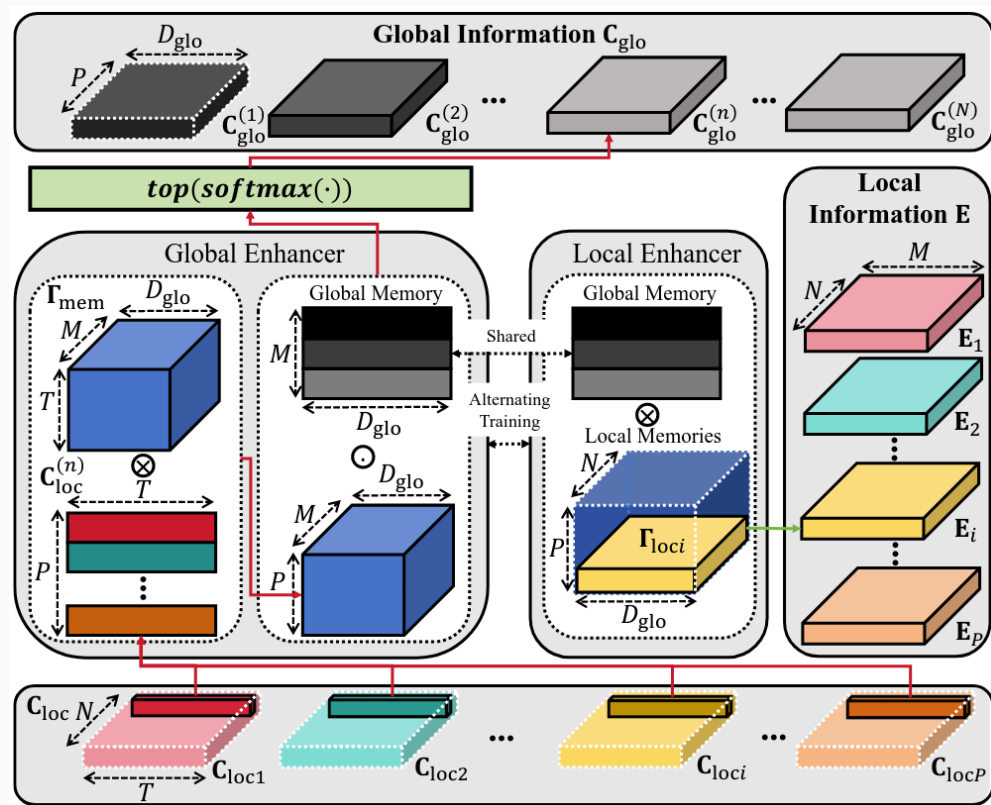
2.4 Alternating Memory Enhancer

Overview

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

Hyperparameters

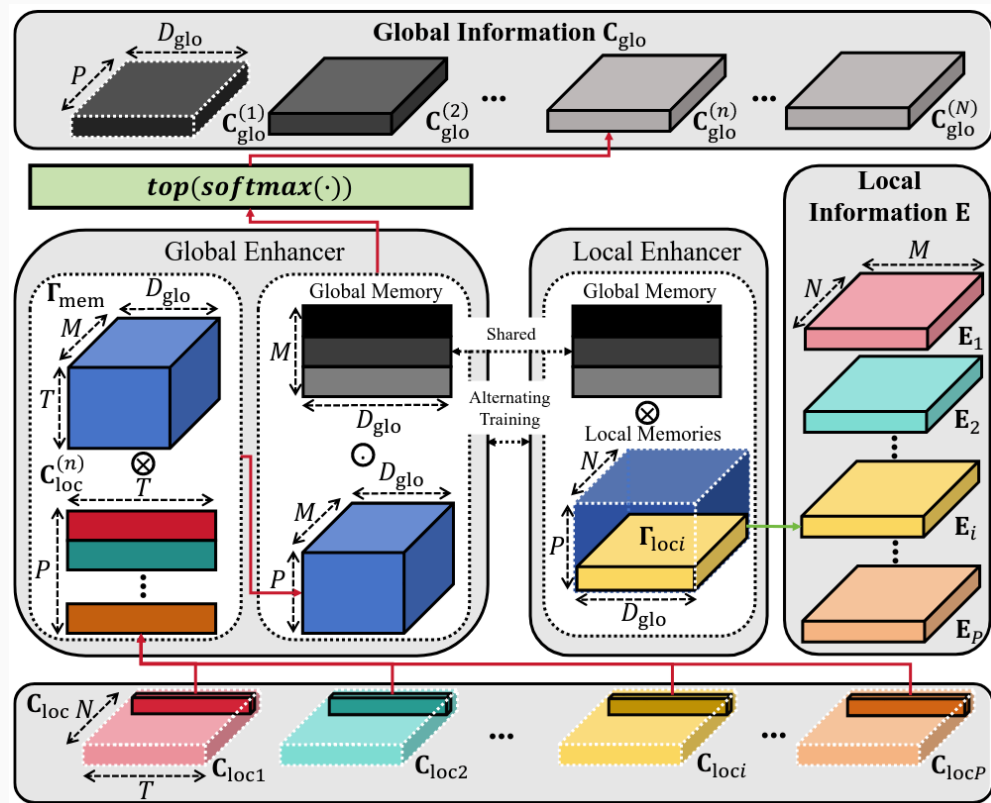
- $M \rightarrow$ number of high level patterns
 - Spikes, seasons, stable
- $D_{\text{glo}} \rightarrow$ richness of patterns



2.4 Alternating Memory Enhancer

Local Enhancer

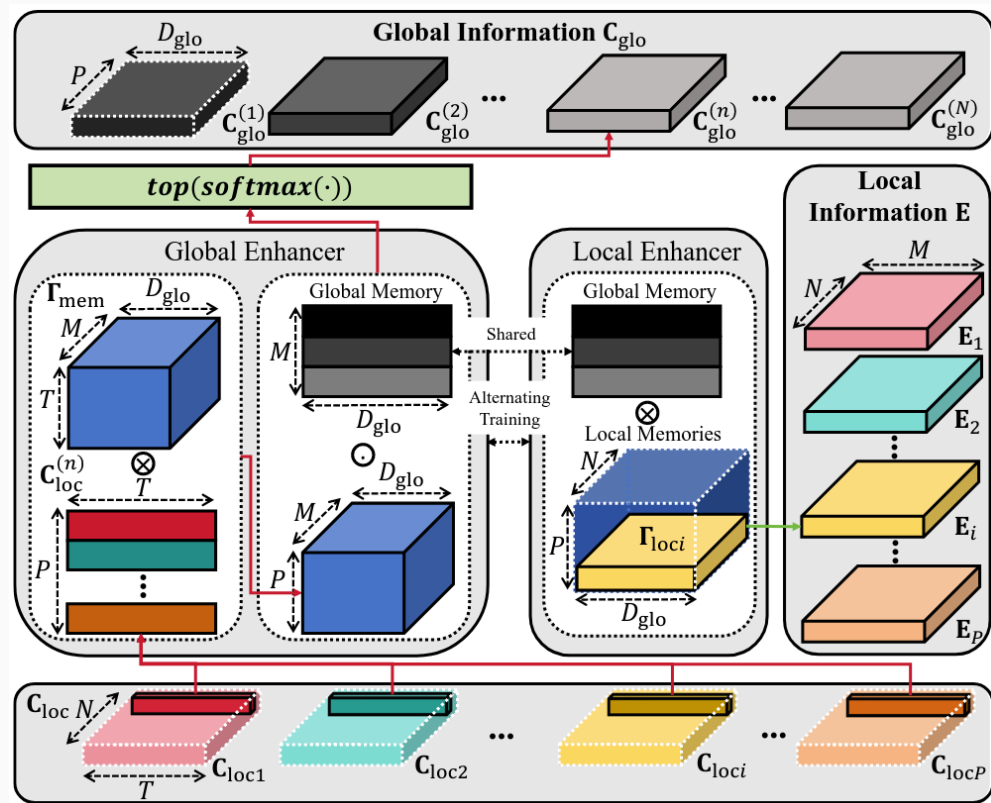
- Local memory regions Γ_{loci}
 - One for each patch
- $P_i \longleftrightarrow \Gamma_{loci} \longrightarrow \Gamma_{loci} \longleftrightarrow E_i$
- Not directly identical
 - E contains global memories influence
 - Defined by C_{loc}
- Memories are **not** information



2.4 Alternating Memory Enhancer

Global Enhancer

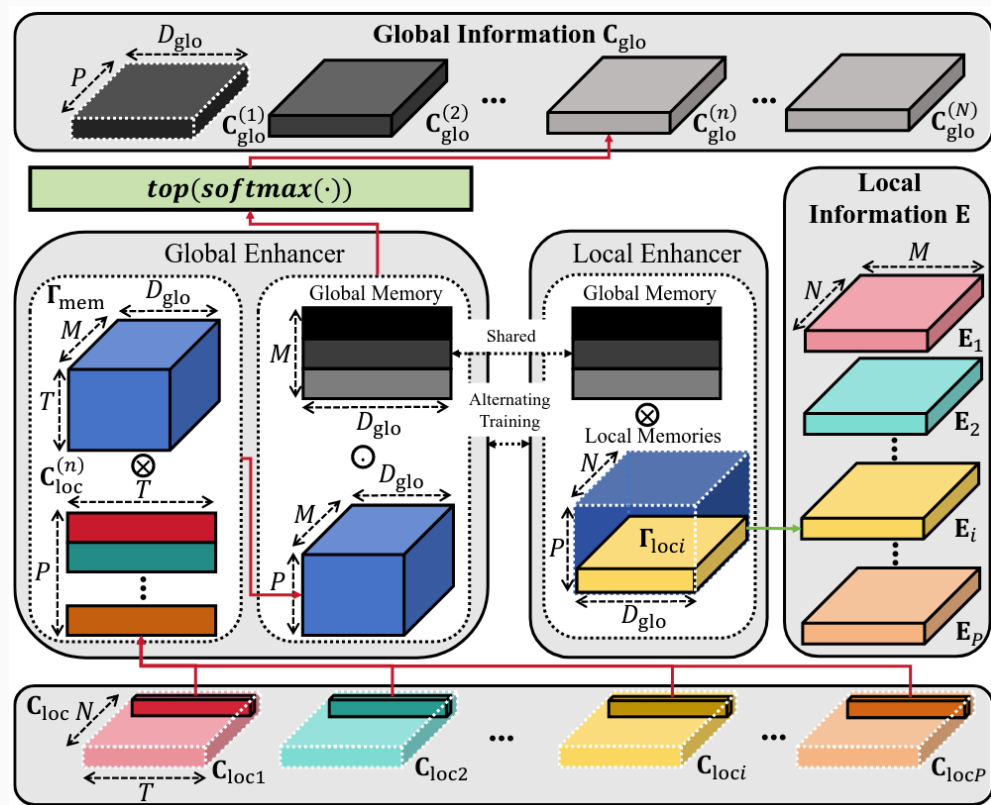
- Learns from locally correlated features
- Γ_{mem} is a large trainable tensor
 - Produces inquiry tensor
 - Recognizes patterns in data
 - The M high level patterns



2.4 Alternating Memory Enhancer

Global Enhancer

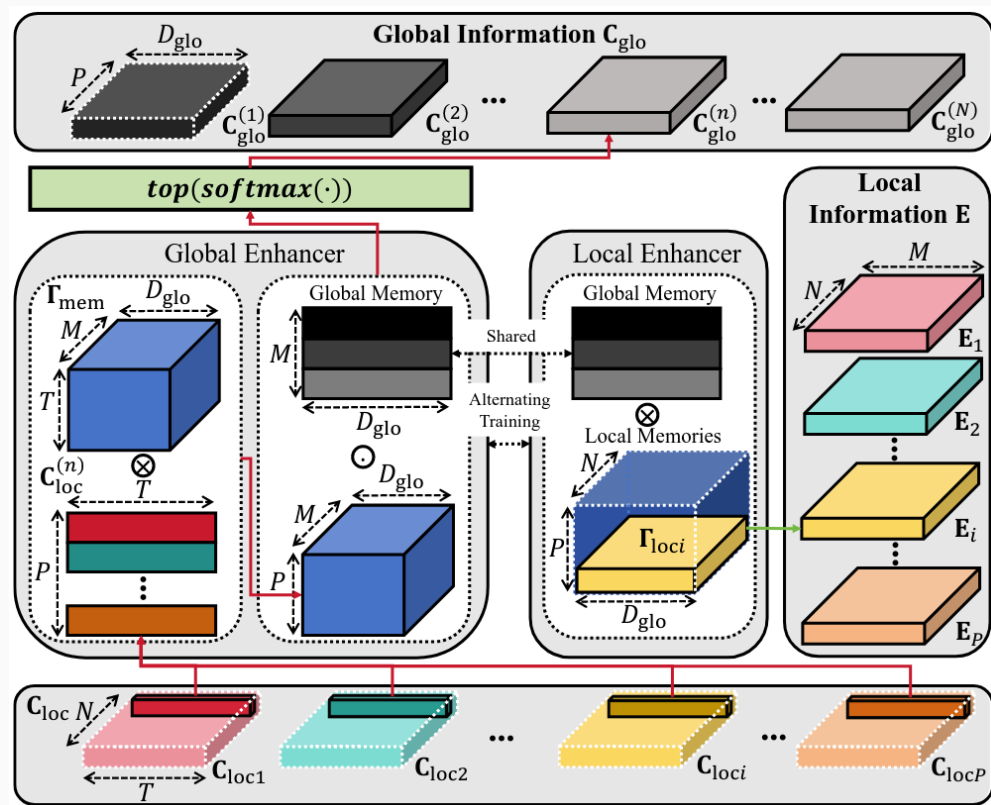
- Learns from locally correlated features
- Γ_{mem} is a large trainable tensor
 - Produces inquiry tensor
 - Recognizes patterns in data
 - The M high level patterns
- Inquiry tensor
 - Prevalence of patterns in local data
 - Similarity scores with global memory



2.4 Alternating Memory Enhancer

Global Enhancer

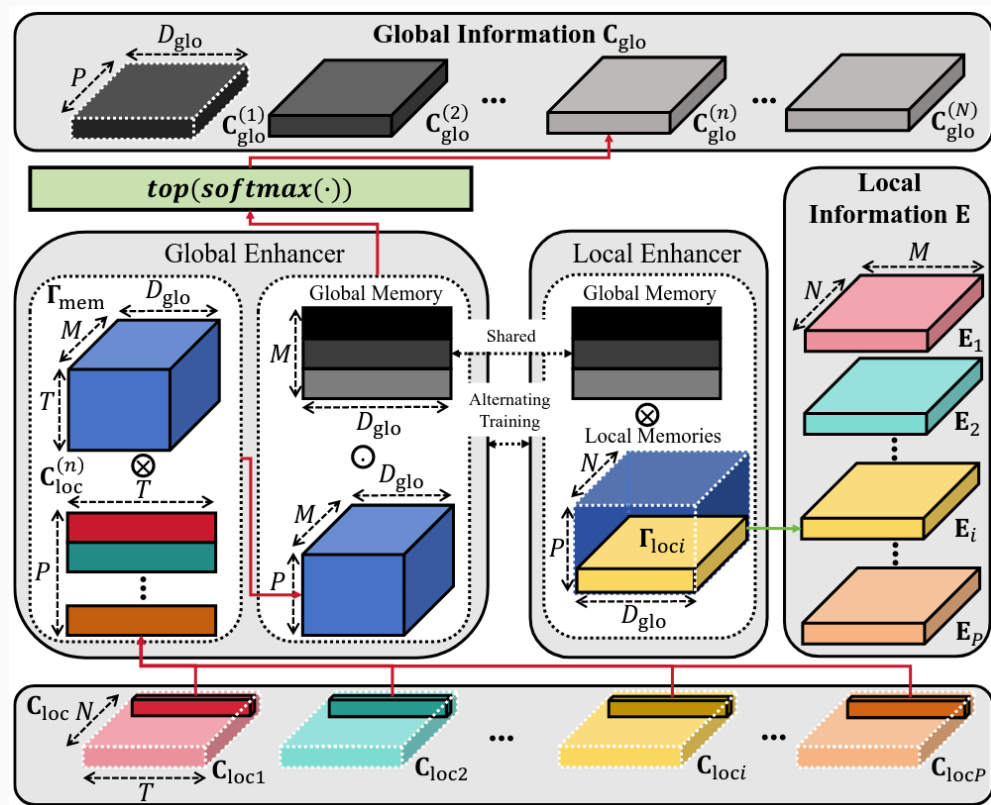
- Learns from locally correlated features
- Γ_{mem} is a large trainable tensor
 - Produces inquiry tensor
 - Recognizes patterns in data
 - The M high level patterns
- Inquiry tensor
 - Prevalence of patterns in local data
 - Similarity scores with global memory
- Probability distribution
 - Importance of pattern



2.4 Alternating Memory Enhancer

Global Enhancer

- Learns from locally correlated features
- Γ_{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_{glo}
 - Scaled based on importance
 - Weighted sum



2.4 Alternating Memory Enhancer

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

2.4 Alternating Memory Enhancer

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 → longer convergence
- Balance convergence
 - Different learning rates
 - LE training more

2.4 Alternating Memory Enhancer

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 ϵ ; 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 Γ_{mem} , and bias \mathbf{b}_{mem}

- 1: *Initialisation:* Initializing local and global memories $\Gamma_{\text{loc}}, \Gamma_{\text{glo}}$, tensor Γ_{mem} , and bias \mathbf{b}_{mem} randomly
- 2: **while** $\Gamma_{\text{loc}}, \Gamma_{\text{glo}}, \Gamma_{\text{mem}}$, and \mathbf{b}_{mem} are not converged **do**
- 3: **for** iteration = 0 to ϵ **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: $\Gamma_{\text{loc}} \leftarrow \Gamma_{\text{loc}} - \eta_{\text{loc}} \nabla_{\Gamma_{\text{loc}}} \mathcal{L}(\hat{\mathbf{F}}, \mathbf{F})$
- 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: $\Gamma_{\text{glo}} \leftarrow \Gamma_{\text{glo}} - \eta_{\text{glo}} \nabla_{\Gamma_{\text{glo}}} \mathcal{L}(\hat{\mathbf{F}}, \mathbf{F})$
- 19: $\Gamma_{\text{mem}} \leftarrow \Gamma_{\text{mem}} - \eta_{\text{glo}} \nabla_{\Gamma_{\text{mem}}} \mathcal{L}(\hat{\mathbf{F}}, \mathbf{F})$,
- 20: $\mathbf{b}_{\text{mem}} \leftarrow \mathbf{b}_{\text{mem}} - \eta_{\text{glo}} \nabla_{\mathbf{b}_{\text{mem}}} \mathcal{L}(\hat{\mathbf{F}}, \mathbf{F})$
- 21: **end while**

3. Experiments

3.1 Noteworthy Details

Datasets

- 7 in total
 - 4 are variants of the same
- 7, 21, 321, and 862 variables
- $H = 336$
- $F = [96, 192, 336, 720]$

3.1 Noteworthy Details

Datasets

- 7 in total
 - 4 are variants of the same
- 7, 21, 321, and 862 variables
- $H = 336$
- $F = [96, 192, 336, 720]$

Comparisons

- Multiple different model architectures
 - Channel independent models
 - Linear models
 - Attention models

3.2 Forecasting Accuracy

Results

- Compare on MSE and MAE
- Bold is best, underline is second best
- Almost always best performance
 - Loses on MSE for low F in one dataset

Models		Memformer		ModernTCN		PatchTST		NLinear		DLinear		iTransformer		CARD		Crossformer		MTGNN	
Metric		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Weather	96	0.151	0.185	0.155	0.201	0.152	0.199	0.182	0.232	0.176	0.237	0.174	0.214	<u>0.150</u>	<u>0.188</u>	0.145	0.211	0.342	0.385
	192	0.197	0.231	0.198	0.245	0.197	0.243	0.225	0.269	0.220	0.282	0.221	0.254	0.202	<u>0.238</u>	0.190	0.259	0.427	0.445
	336	0.247	0.274	0.251	0.286	<u>0.249</u>	0.283	0.271	0.301	0.265	0.319	0.278	0.296	0.260	<u>0.282</u>	0.259	0.326	0.506	0.523
	720	0.318	0.326	0.321	0.336	0.320	0.335	0.338	0.348	0.323	0.362	0.358	0.347	0.343	0.353	0.332	0.382	0.510	0.527
Traffic	96	0.361	0.230	0.368	0.253	<u>0.367</u>	0.251	0.410	0.279	0.410	0.282	0.395	0.268	0.419	0.269	0.511	0.292	0.516	0.308
	192	0.381	0.239	0.384	0.261	0.385	<u>0.259</u>	0.423	0.284	0.423	0.287	0.417	0.276	0.443	0.276	0.523	0.311	0.534	0.324
	336	0.394	0.245	<u>0.397</u>	0.270	0.398	<u>0.265</u>	0.435	0.290	0.436	0.296	0.433	0.283	0.460	0.283	0.530	0.300	0.540	0.335
	720	0.432	0.267	0.440	0.296	<u>0.434</u>	<u>0.287</u>	0.464	0.307	0.466	0.315	0.467	0.302	0.490	0.299	0.573	0.313	0.557	0.343
Electricity	96	0.130	0.217	<u>0.131</u>	0.228	0.130	0.222	0.141	0.237	0.140	0.237	0.132	0.228	0.141	0.233	0.186	0.281	0.202	0.314
	192	0.147	0.232	0.150	0.242	<u>0.148</u>	<u>0.240</u>	0.154	0.248	0.153	0.249	0.154	0.249	0.160	0.250	0.208	0.300	0.266	0.349
	336	0.162	0.249	0.171	0.265	<u>0.167</u>	0.261	0.171	0.265	0.169	0.267	0.172	0.267	0.173	0.263	0.323	0.369	0.328	0.373
	720	0.199	0.281	0.203	0.294	<u>0.202</u>	0.291	0.210	0.297	0.203	0.301	0.204	0.296	0.197	0.284	0.404	0.423	0.422	0.410
ETTh1	96	0.362	0.385	0.382	0.401	0.375	0.399	0.374	0.394	0.375	0.399	0.386	0.405	0.383	0.391	0.377	0.419	0.401	0.442
	192	0.386	0.404	0.420	0.424	0.414	0.421	0.408	0.415	0.405	0.416	0.441	0.436	0.435	0.420	0.410	0.439	0.587	0.601
	336	0.402	0.421	0.427	0.434	0.431	0.436	0.429	<u>0.427</u>	0.439	0.443	0.487	0.458	0.479	0.442	0.440	0.461	0.736	0.643
	720	0.436	0.452	0.450	0.461	0.449	0.466	0.440	<u>0.453</u>	0.472	0.490	0.503	0.491	0.471	0.461	0.519	0.524	0.916	0.750
ETTh2	96	0.264	0.321	0.276	0.342	0.274	0.336	0.277	0.338	0.289	0.353	0.297	0.349	0.281	0.330	0.770	0.529	0.735	0.643
	192	0.314	0.358	0.340	0.381	0.339	0.379	0.344	0.381	0.383	0.418	0.380	0.400	0.363	0.381	0.848	0.657	0.859	0.717
	336	0.312	0.364	0.329	0.378	0.331	0.380	0.357	0.400	0.448	0.465	0.428	0.432	0.411	0.418	0.859	0.674	1.050	0.849
	720	0.374	0.410	0.392	0.433	0.379	0.422	0.394	0.436	0.605	0.551	0.427	0.445	0.416	0.431	1.221	0.825	1.336	0.963
ETTm1	96	0.285	0.336	0.292	0.346	0.290	0.342	0.306	0.348	0.299	0.343	0.334	0.368	0.316	0.347	0.320	0.373	0.428	0.446
	192	0.323	0.358	0.332	0.368	<u>0.332</u>	0.369	0.349	0.375	0.335	0.365	0.377	0.391	0.363	0.370	0.372	0.411	0.551	0.505
	336	0.365	0.381	0.367	0.393	0.366	0.392	0.375	0.388	0.369	<u>0.386</u>	0.426	0.420	0.392	0.390	0.429	0.441	0.706	0.622
	720	0.419	0.409	0.422	0.429	0.420	0.424	0.433	0.422	0.425	<u>0.421</u>	0.491	0.459	0.458	0.425	0.573	0.531	0.982	0.764
ETTm2	96	0.160	0.245	0.166	0.256	0.165	0.255	0.167	0.255	0.167	0.260	0.180	0.264	0.169	0.248	0.254	0.348	0.442	0.483
	192	0.215	0.285	0.222	0.293	0.220	0.293	0.221	0.293	0.224	0.303	0.250	0.309	0.234	0.292	0.370	0.433	0.642	0.570
	336	0.263	0.317	0.276	0.327	0.278	0.329	0.274	0.327	0.281	0.342	0.311	0.348	0.294	0.339	0.511	0.527	0.726	0.658
	720	0.350	0.372	0.365	0.383	0.367	0.385	0.368	0.384	0.397	0.421	0.412	0.407	0.390	0.388	0.901	0.689	1.139	0.862

3.3 Ablation Study

Overview

- Are components contributing?
- Experiment without
 - Graph learning
 - GRU
 - Local
 - Global
 - Sharing
 - Alternating

3.3 Ablation Study

Overview

- Are components contributing?
- Experiment without
 - Graph learning
 - GRU
 - Local
 - Global
 - Sharing
 - Alternating

Models		Memformer		w/o Graph		w/o Recurrent		w/o Local		w/o Global		w/o Sharing		w/o Alternating	
Metric		MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	MAE
Weather	96	0.151	0.185	0.152	0.197	0.155	0.194	0.159	0.204	0.153	0.195	0.151	0.187	0.155	0.199
	192	0.197	0.231	0.200	0.235	0.204	0.248	0.199	0.255	0.202	0.238	0.199	0.234	0.202	0.235
	336	0.247	0.274	0.252	0.279	0.254	0.298	0.257	0.307	0.251	0.284	0.252	0.279	0.255	0.286
	720	0.318	0.326	0.334	0.341	0.332	0.355	0.364	0.380	0.323	0.333	0.324	0.335	0.334	0.360
Electricity	96	0.130	0.217	0.133	0.226	0.132	0.224	0.132	0.243	0.131	0.223	0.131	0.223	0.138	0.235
	192	0.147	0.232	0.154	0.245	0.155	0.248	0.153	0.250	0.150	0.238	0.152	0.241	0.158	0.252
	336	0.162	0.249	0.169	0.260	0.174	0.268	0.179	0.270	0.169	0.258	0.170	0.261	0.181	0.274
	720	0.199	0.281	0.208	0.299	0.220	0.317	0.231	0.341	0.205	0.293	0.210	0.300	0.229	0.339
ETTh2	96	0.264	0.321	0.271	0.329	0.269	0.326	0.322	0.369	0.266	0.324	0.266	0.324	0.294	0.347
	192	0.314	0.358	0.328	0.365	0.325	0.362	0.458	0.478	0.320	0.364	0.318	0.361	0.372	0.401
	336	0.312	0.364	0.329	0.376	0.334	0.381	0.530	0.517	0.317	0.370	0.319	0.370	0.380	0.419
	720	0.374	0.410	0.379	0.421	0.401	0.437	0.705	0.627	0.385	0.422	0.388	0.425	0.437	0.463

Results

- All component are contributing

3.4 Challenge 1

Disrupted Correlations

- Robustness
- Introduce outliers
 - Different amounts
 - Independent
 - Dependent

3.4 Challenge 1

Disrupted Correlations

- Robustness
- Introduce outliers
 - Different amounts
 - Independent
 - Dependent
- Results
 - Performed the best
 - Mitigate both types of outliers

3.4 Challenge 1

Disrupted Correlations

- Robustness
- Introduce outliers
 - Different amounts
 - Independent
 - Dependent
- Results
 - Performed the best
 - Mitigate both types of outliers

Dynamic Correlations

- Introduce dynamic correlations
 - Different amounts

3.4 Challenge 1

Disrupted Correlations

- Robustness
- Introduce outliers
 - Different amounts
 - Independent
 - Dependent
- Results
 - Performed the best
 - Mitigate both types of outliers

Dynamic Correlations

- Introduce dynamic correlations
 - Different amounts
- Results
 - Performed the best

4. Critique

4.1 Preprocessing

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

μ is the mean

σ is the variance

c ensures numerical stability

4.1 Preprocessing

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

μ is the mean

σ 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, μ and σ denote the mean and variance of the sample, respectively, and “constant” is a small positive real number included to ensure numerical stability.

4.1 Preprocessing

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

μ is the mean

σ 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, μ and σ 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?
 - σ is conventional notation for standard deviation
 - σ^2 is conventional notation for variance

4.1 Preprocessing

What is going on?

4.1 Preprocessing

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`

4.1 Preprocessing

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`

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


4.1 Preprocessing

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`

```
from sklearn.preprocessing import StandardScaler
class ...:
    def __read_data__(self):
        self.scaler = 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

4.1 Preprocessing

What are they actually doing?

Preprocessing

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

H is the historical horizon

μ is the mean

σ is the variance

c ensures numerical stability

StandardScaler

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

x is the sample

μ is the mean

σ is the standard deviation

4.1 Preprocessing

What are they actually doing?

Preprocessing

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

H is the historical horizon

μ is the mean

σ 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

μ is the mean

σ is the standard deviation

4.1 Preprocessing

What are they actually doing?

Preprocessing

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

H is the historical horizon

μ is the mean

σ is the variance

c ensures numerical stability

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

StandardScaler

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

x is the sample

μ is the mean

σ is the standard deviation

4.1 Preprocessing

What are they actually doing?

Preprocessing

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

H is the historical horizon

μ is the mean

σ is the variance

c ensures numerical stability

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

StandardScaler

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

x is the sample

μ is the mean

σ is the standard deviation

4.1 Preprocessing

What are they actually doing?

Preprocessing

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

H is the historical horizon

μ is the mean

σ is the variance

c ensures numerical stability

StandardScaler

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

x is the sample

μ is the mean

σ 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

4.2 Notation

Inconsistencies

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

4.2 Notation

Inconsistencies

- C_{glo} is global memory
- C_{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

4.2 Notation

Inconsistencies

- C_{glo} is global memory
- C_{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.

4.2 Notation

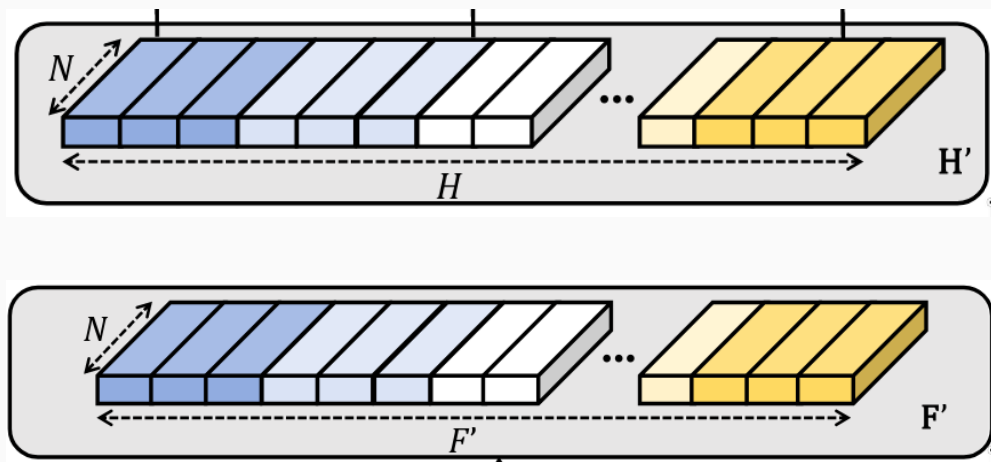
Inconsistencies

- C_{glo} is global memory
- C_{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.

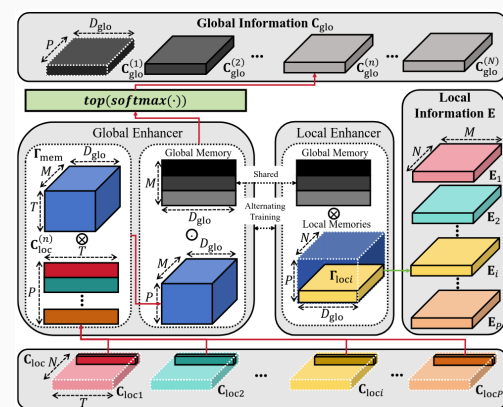
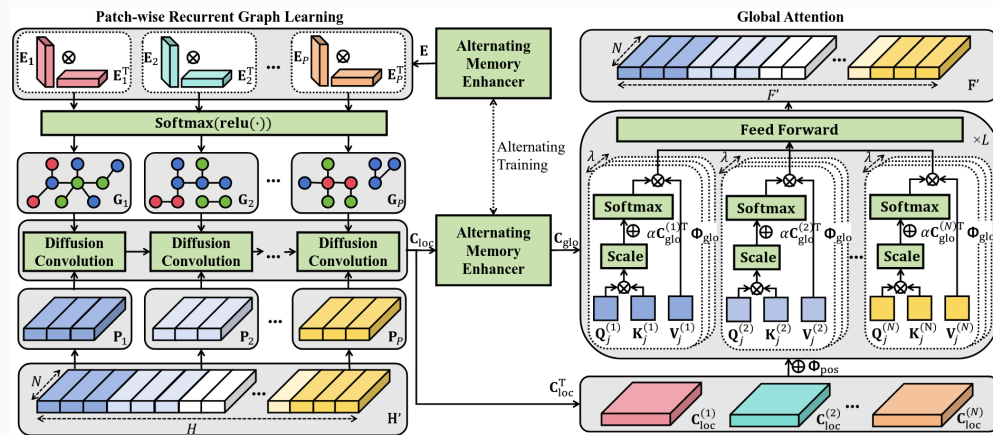


5. Praise

5.1 Figures

Colors

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



5.1 Figures

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

