

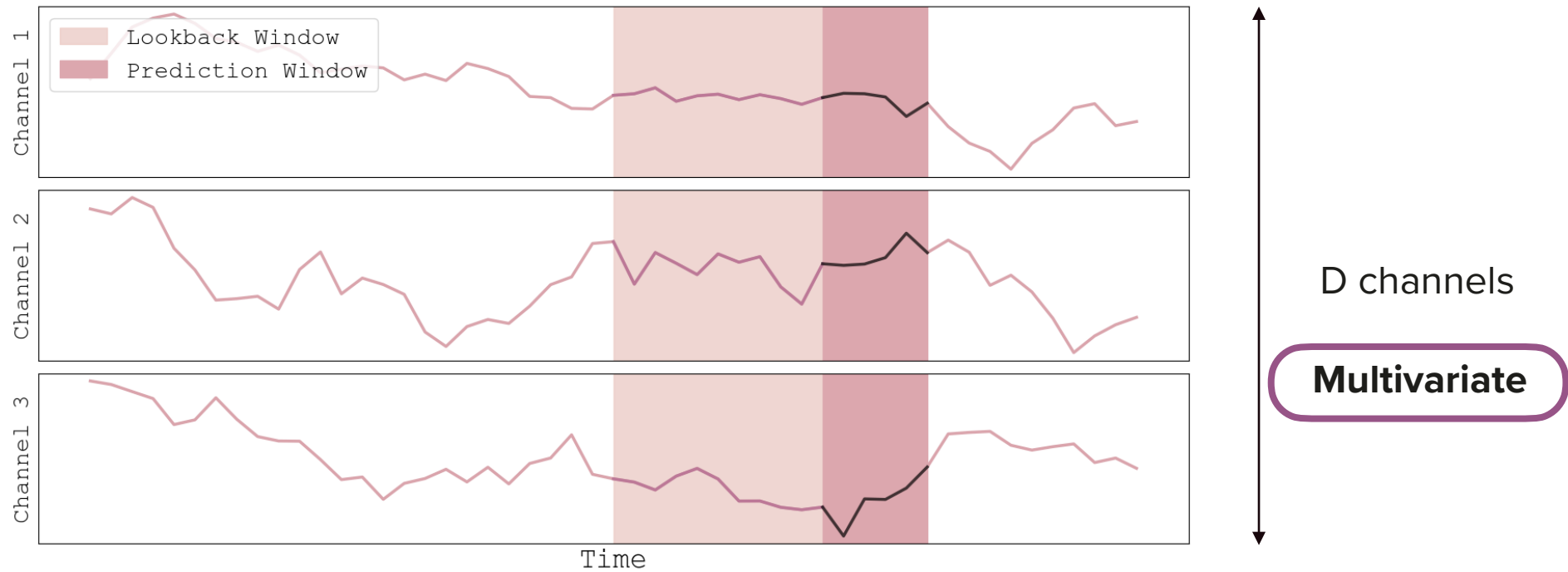
SAMformer : Unlocking the Potential of Transformers in Time-Series Forecasting

Romain ILBERT, Ambroise Odonnat, Vasilii Feofanov, Aladin Virmaux,
Giuseppe Paolo, Ievgen Redko, Themis Palpanas



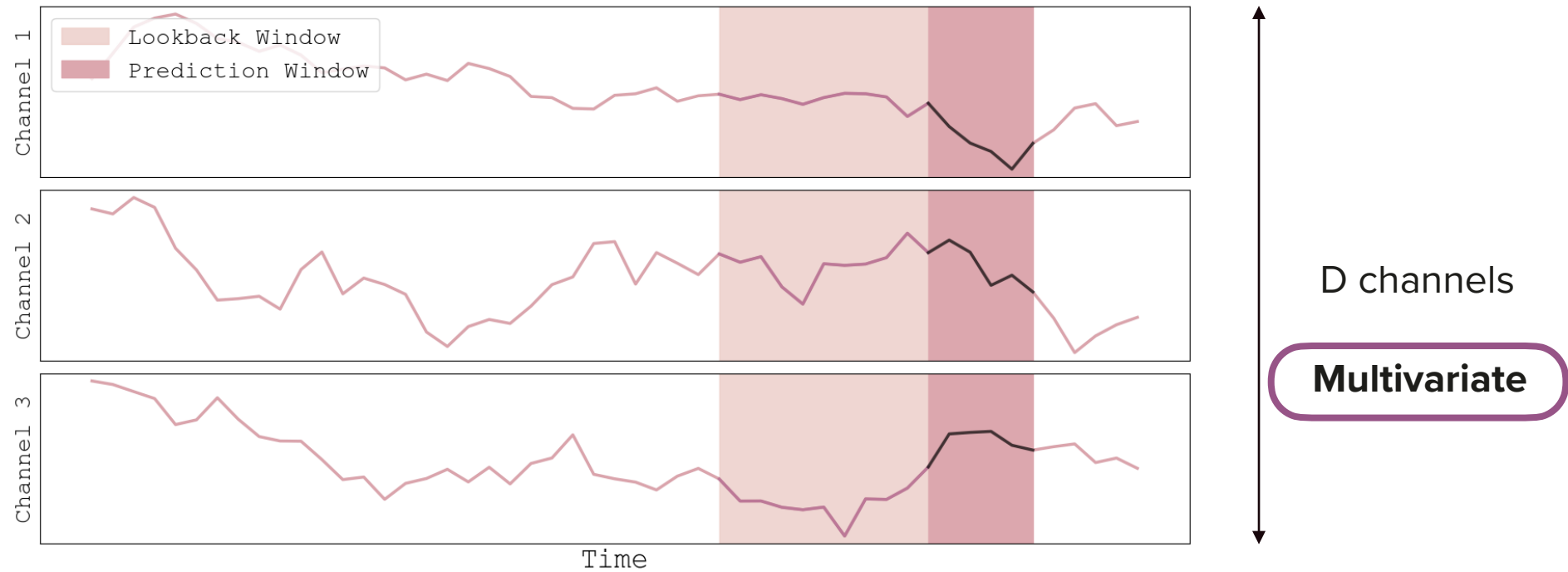
Time Series Forecasting : Problem Setup

Multivariate Time Series Forecasting



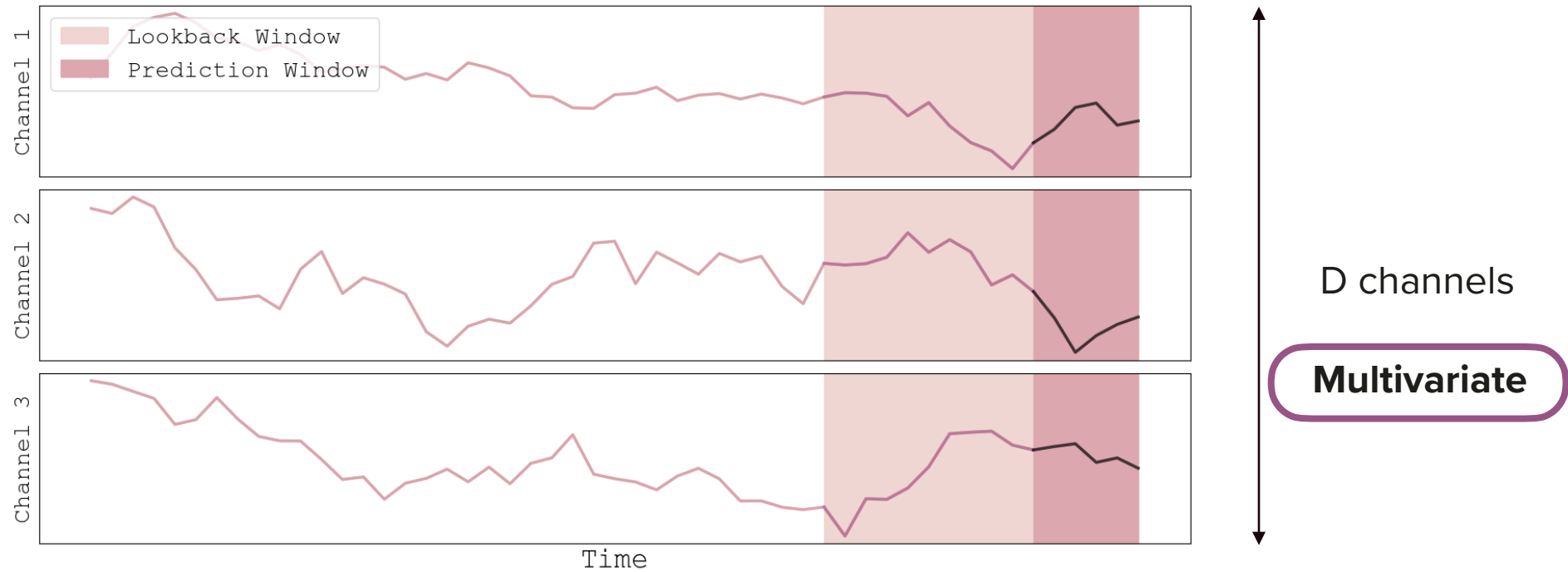
Time Series Forecasting : Problem Setup

Multivariate Time Series Forecasting



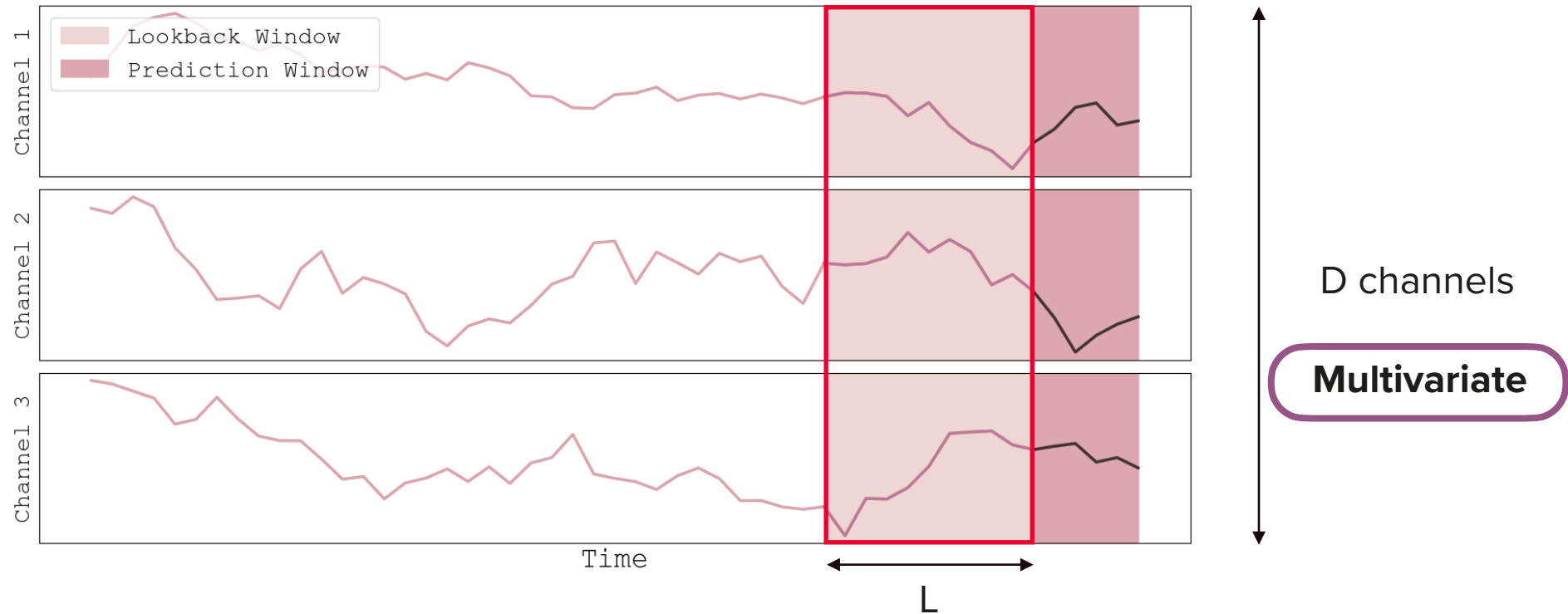
Time Series Forecasting : Problem Setup

Multivariate Time Series Forecasting



Time Series Forecasting : Problem Setup

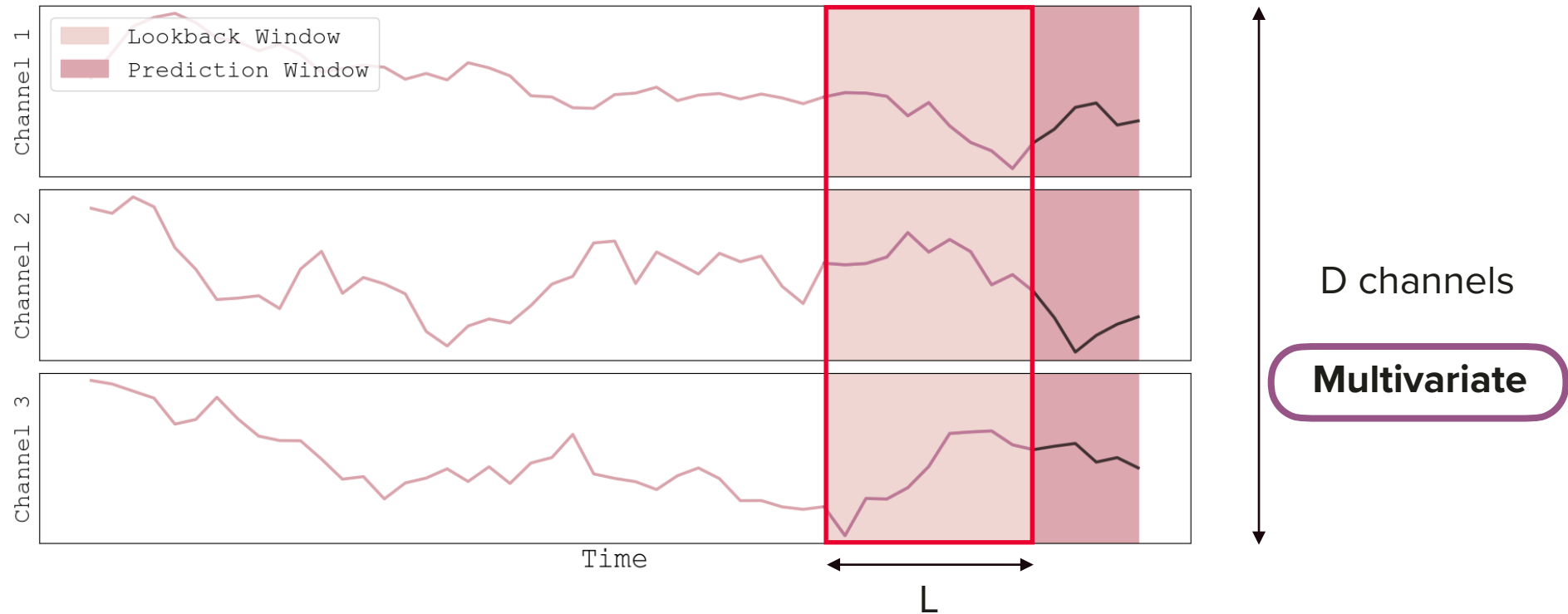
Multivariate Time Series Forecasting



Inputs : $X \in \mathbb{R}^{D \times L}$

Time Series Forecasting : Problem Setup

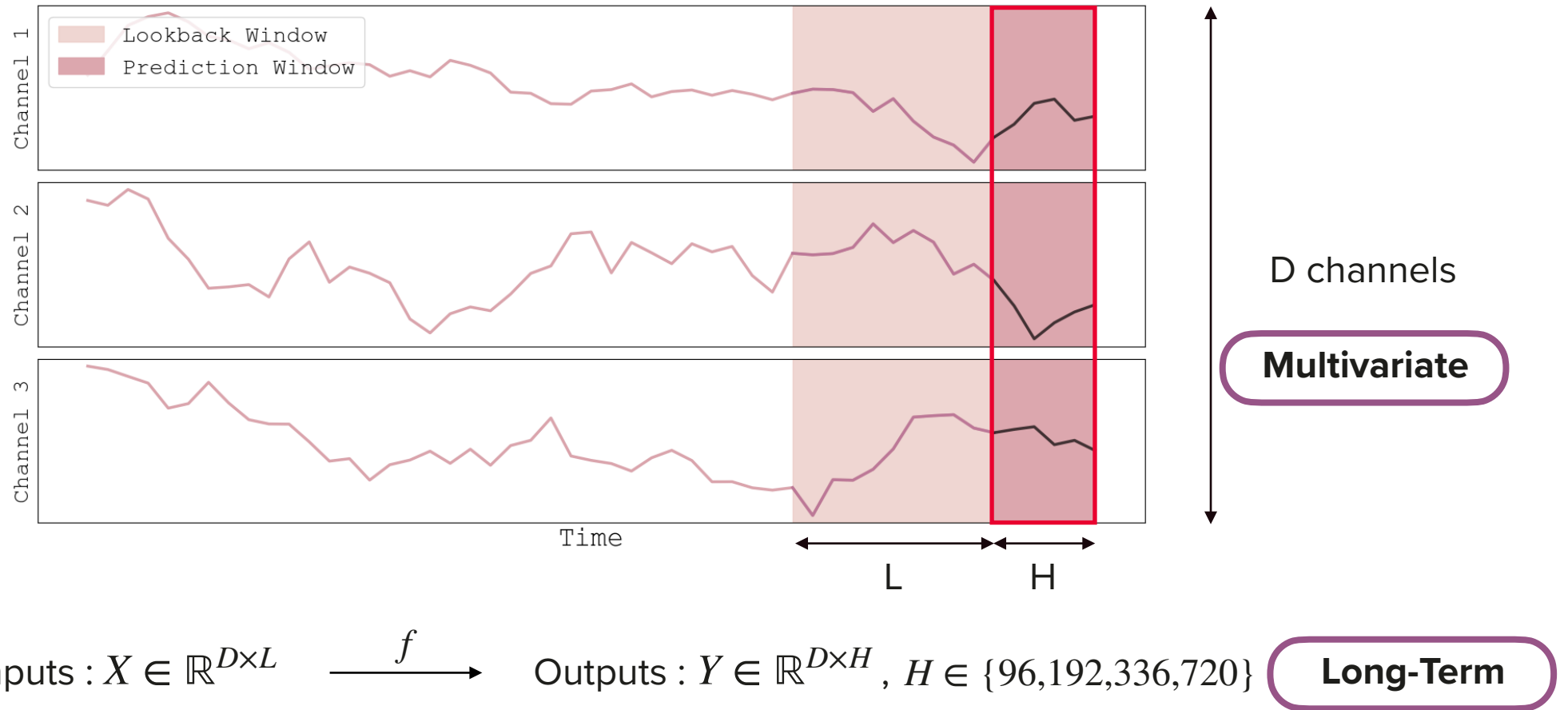
Multivariate Time Series Forecasting



$$\text{Inputs : } X \in \mathbb{R}^{D \times L} \xrightarrow{f}$$

Time Series Forecasting : Problem Setup

Multivariate Time Series Forecasting



Failure of Transformers in Time Series Forecasting

Main conclusions from *Are Transformers Effective for Time Series Forecasting** ?

1. Transformer-based methods don't work well in forecasting

LogSparse and convolutional self-attention @**LogTrans**

ProbSparse and distilling self-attention @**Informer**

Series auto-correlation with decomposition @**Autoformer**

Multi-resolution pyramidal attention @**Pyraformer**

Frequency enhanced block with decomposition @**FEDformer**

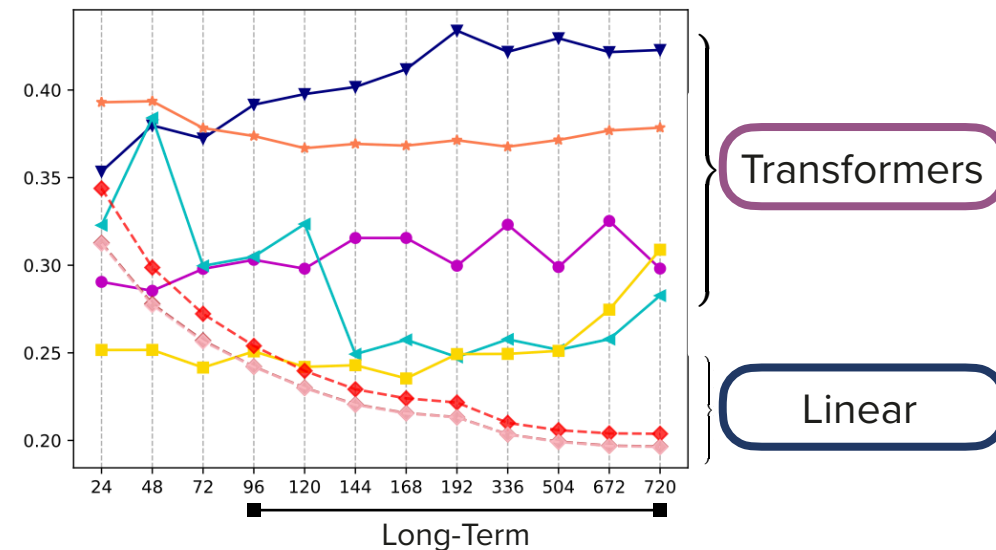
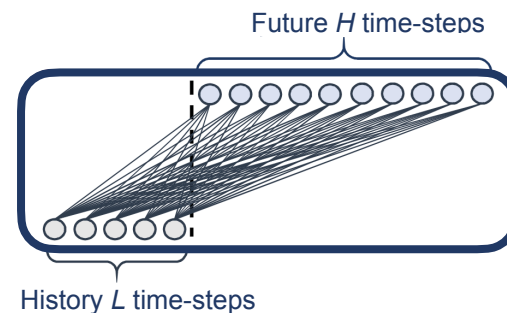
Failure of Transformers in Time Series Forecasting

Main conclusions from *Are Transformers Effective for Time Series Forecasting** ?

1. Transformer-based methods don't work well in forecasting
2. A simple Linear Model surpasses all of them in long-term forecasting

- LogSparse and convolutional self-attention @**LogTrans**
- ProbSparse and distilling self-attention @**Informer**
- Series auto-correlation with decomposition @**Autoformer**
- Multi-resolution pyramidal attention @**Pyraformer**
- Frequency enhanced block with decomposition @**FEDformer**

VS.



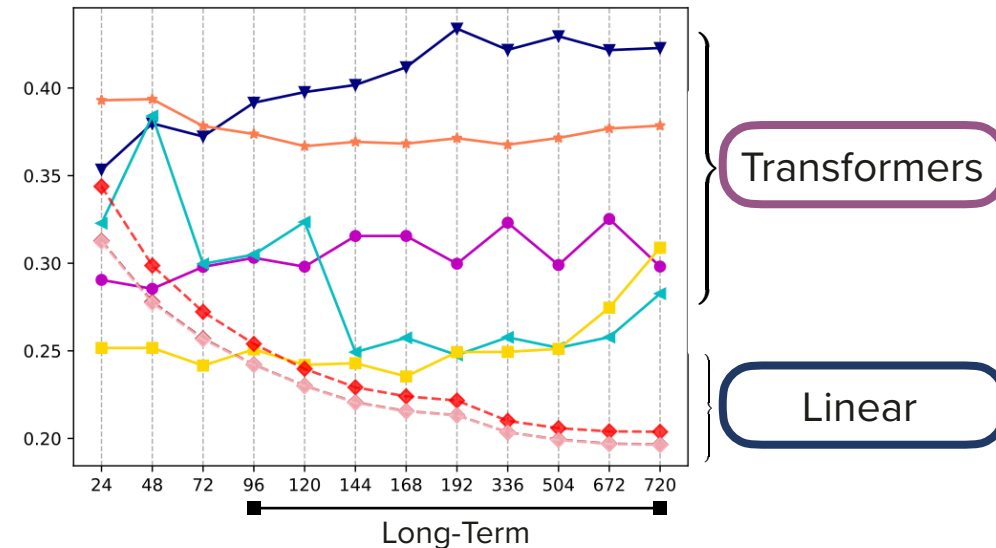
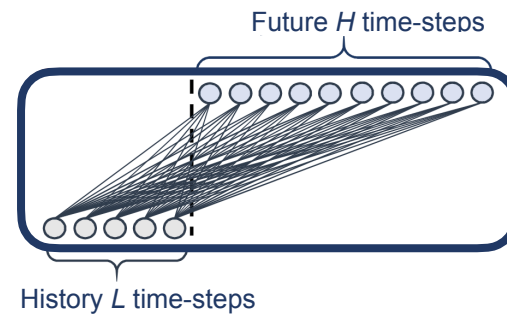
Failure of Transformers in Time Series Forecasting

Main conclusions from *Are Transformers Effective for Time Series Forecasting** ?

1. Transformer-based methods don't work well in forecasting
2. A simple Linear Model surpasses all of them in long-term forecasting

- LogSparse and convolutional self-attention @**LogTrans**
- ProbSparse and distilling self-attention @**Informer**
- Series auto-correlation with decomposition @**Autoformer**
- Multi-resolution pyramidal attention @**Pyraformer**
- Frequency enhanced block with decomposition @**FEDformer**

VS.

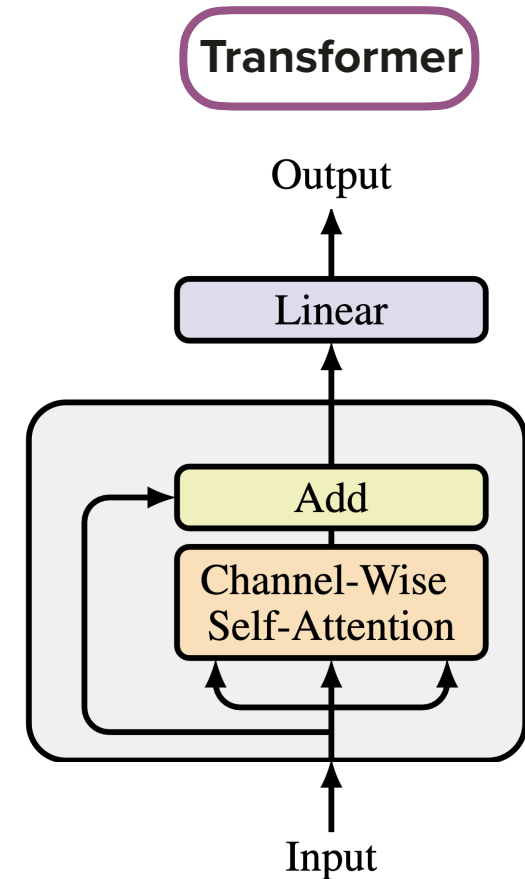


Yet Transformers dominate in NLP and vision... Why?

Failure of Transformers : A simple forecasting example

Context :

- Consider a simple forecasting problem $\mathbf{Y} = \mathbf{XW}_{\text{toy}} + \epsilon$
with $L=512$, $H=96$ and $D=7$



Failure of Transformers : A simple forecasting example

Context :

- Consider a simple forecasting problem $\mathbf{Y} = \mathbf{X}\mathbf{W}_{\text{toy}} + \epsilon$
with $L=512$, $H=96$ and $D=7$

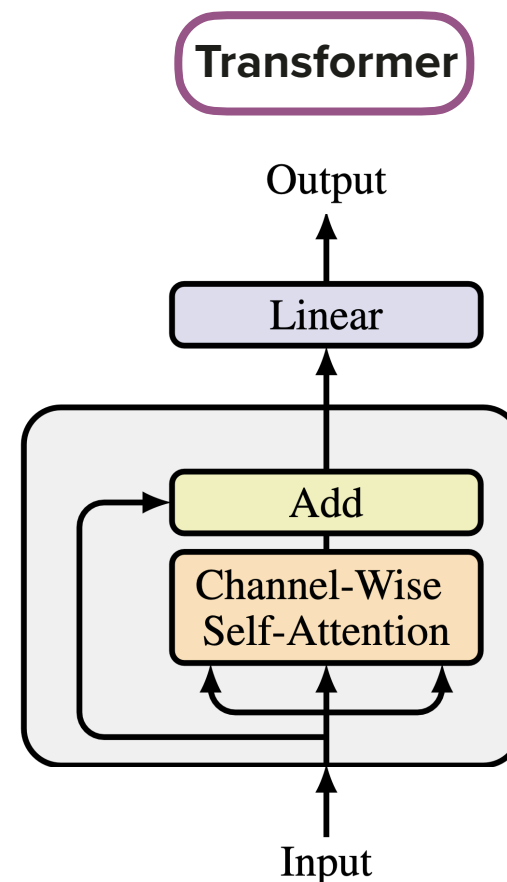
Why channel-wise attention ? :

$$f(\mathbf{X}) = [\mathbf{X} + \text{attention}(\mathbf{X})]\mathbf{W}$$

$D \times L$ $D \times L$ $L \times H$

$$\text{attention}(\mathbf{X}) = \mathbf{A}(\mathbf{X})\mathbf{X}\mathbf{W}_V\mathbf{W}_O$$

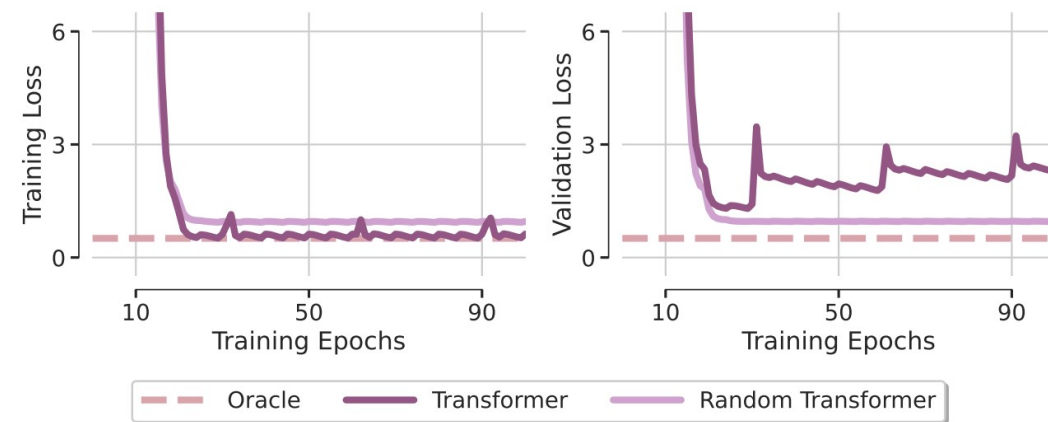
$$\mathbf{A}(\mathbf{X}) \in \mathbb{R}^{D \times D}$$



Failure of Transformers : A simple forecasting example

Conclusions :

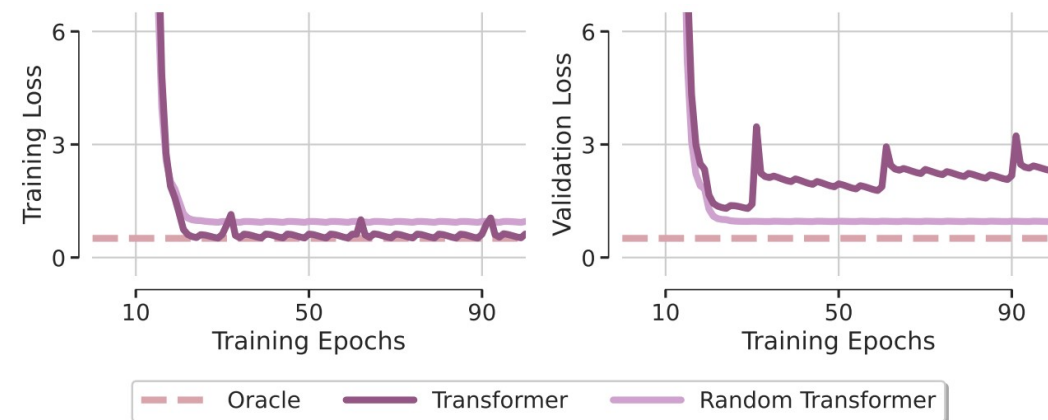
1. Our Transformer severely overfits...
2. ... and works better if we freeze the attention...
3. ... because the attention get stuck in a bad local minimum and does not move afterwards



Failure of Transformers : A simple forecasting example

Conclusions :

1. Our Transformer severely overfits...
2. ... and works better if we freeze the attention...
3. ... because the attention get stuck in a bad local minimum and does not move afterwards

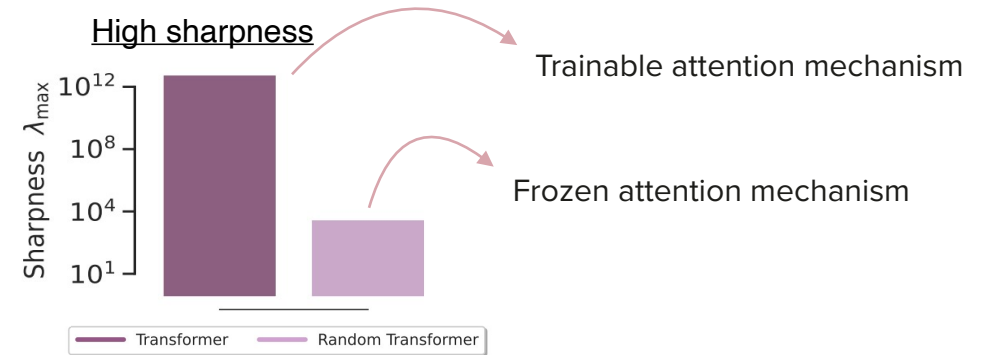
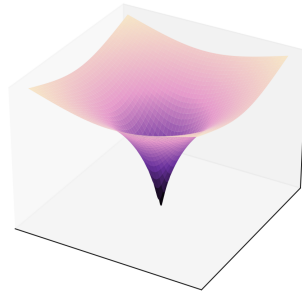


Pathological behavior suggesting sharp local minima

Failure of Transformers : A simple forecasting example

Why transformers fail?

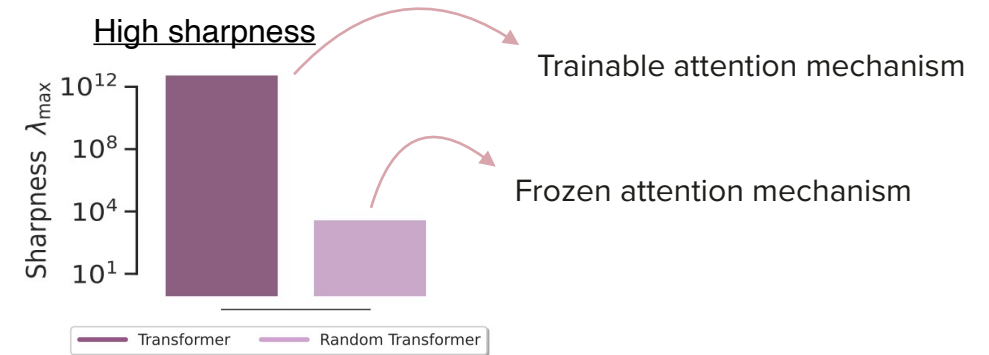
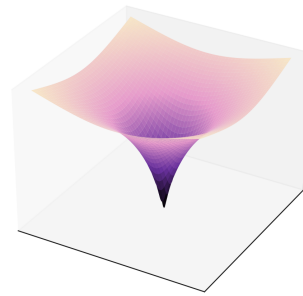
1. Transformers have a sharp loss landscape



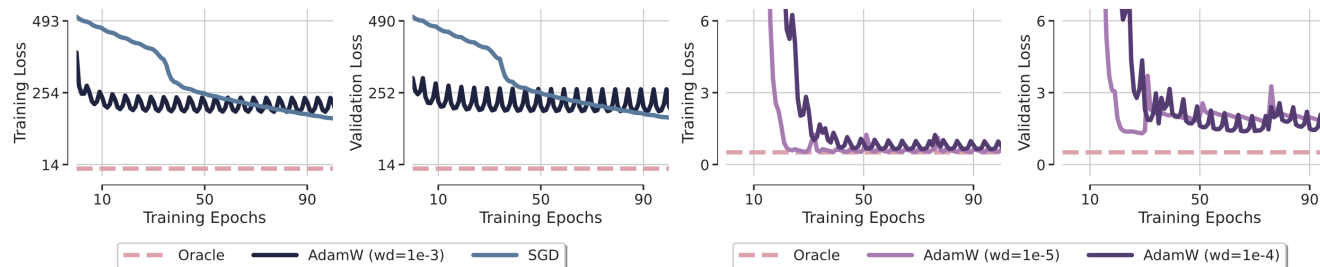
Failure of Transformers : A simple forecasting example

Why transformers fail?

1. Transformers have a sharp loss landscape



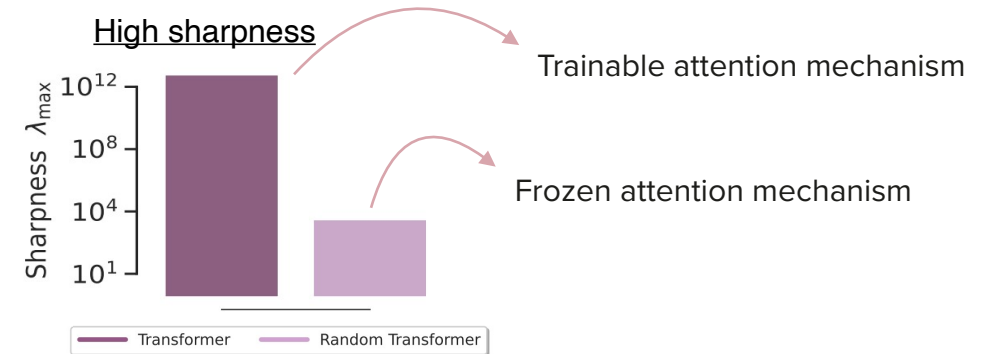
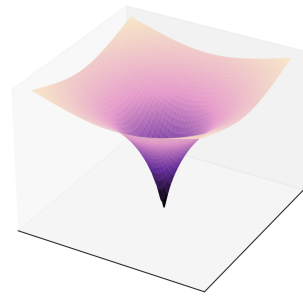
2. And no change in the optimizer helps to solve this



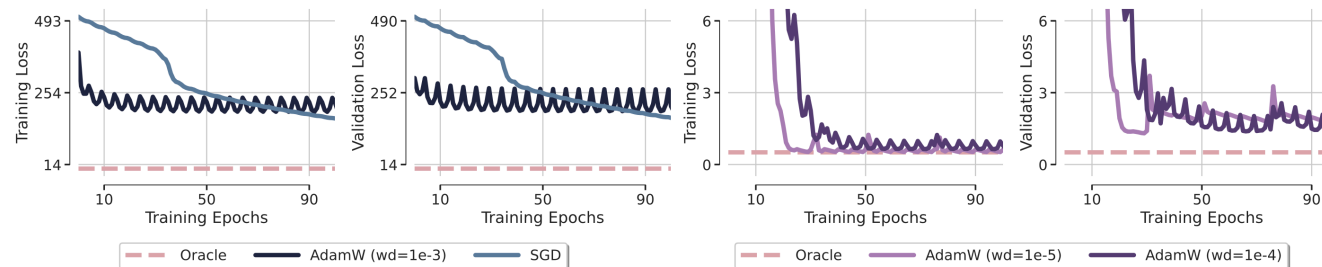
Failure of Transformers : A simple forecasting example

Why transformers fail?

1. Transformers have a sharp loss landscape



2. And no change in the optimizer helps to solve this



3. Well-known in NLP and vision (Chen et al., 2022, Zhai et al. 2023), ignored in TS Forecasting

Failure of Transformers : Solution

Sharpness - Aware Minimization (Foret et al. 2021, Chen et al. 2022)

- Smooths the loss landscape => flatter, more generalizable local minima

$$\mathcal{L}_{\text{train}}^{\text{SAM}}(\omega) = \max_{\|\epsilon\| < \rho} \mathcal{L}_{\text{train}}(\omega + \epsilon)$$

Failure of Transformers : Solution

Sharpness - Aware Minimization (Foret et al. 2021, Chen et al. 2022)

- Smooths the loss landscape => flatter, more generalizable local minima

$$\mathcal{L}_{\text{train}}^{\text{SAM}}(\omega) = \max_{\|\epsilon\| < \rho} \mathcal{L}_{\text{train}}(\omega + \epsilon)$$

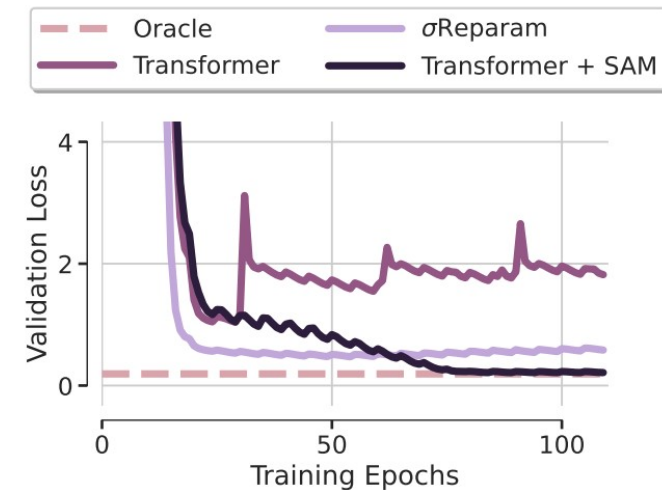
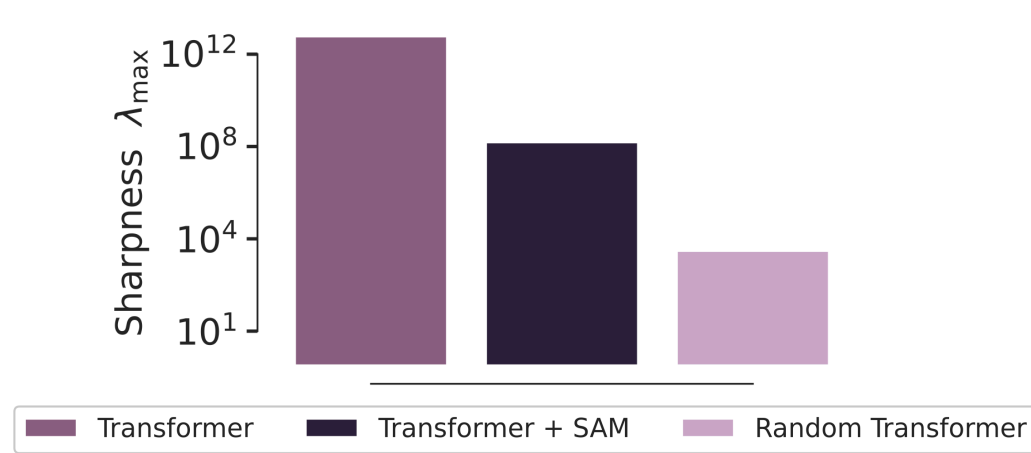


Failure of Transformers : Solution

Sharpness - Aware Minimization (Foret et al. 2021, Chen et al. 2022)

- Smooths the loss landscape => flatter, more generalizable local minima

$$\mathcal{L}_{\text{train}}^{\text{SAM}}(\omega) = \max_{\|\epsilon\| < \rho} \mathcal{L}_{\text{train}}(\omega + \epsilon)$$

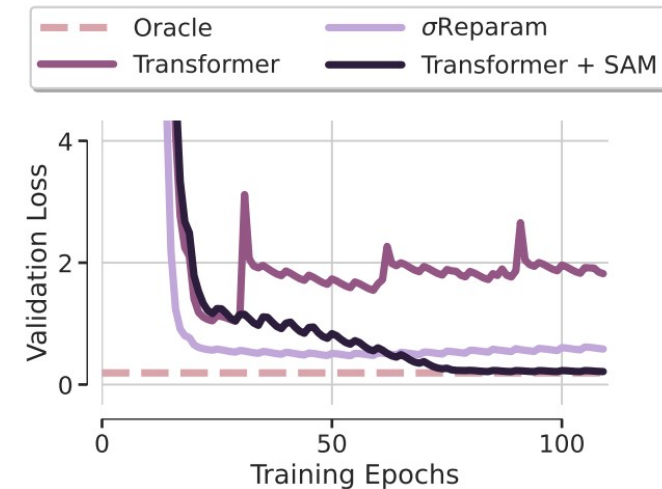
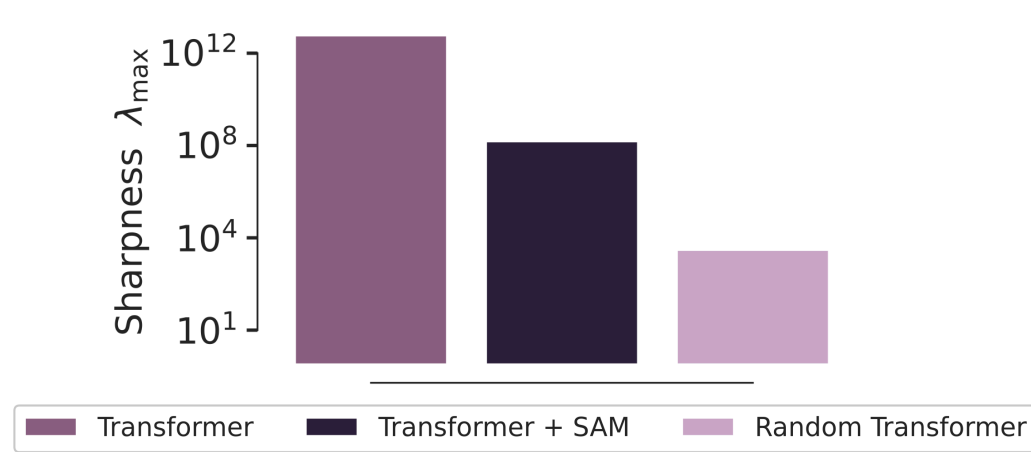


Failure of Transformers : Solution

Sharpness - Aware Minimization (Foret et al. 2021, Chen et al. 2022)

- Smooths the loss landscape => flatter, more generalizable local minima

$$\mathcal{L}_{\text{train}}^{\text{SAM}}(\omega) = \max_{\|\epsilon\| < \rho} \mathcal{L}_{\text{train}}(\omega + \epsilon)$$



SAM = Desired Solution

**Congrats, you now know how to solve
a linear forecasting problem with transformers!**



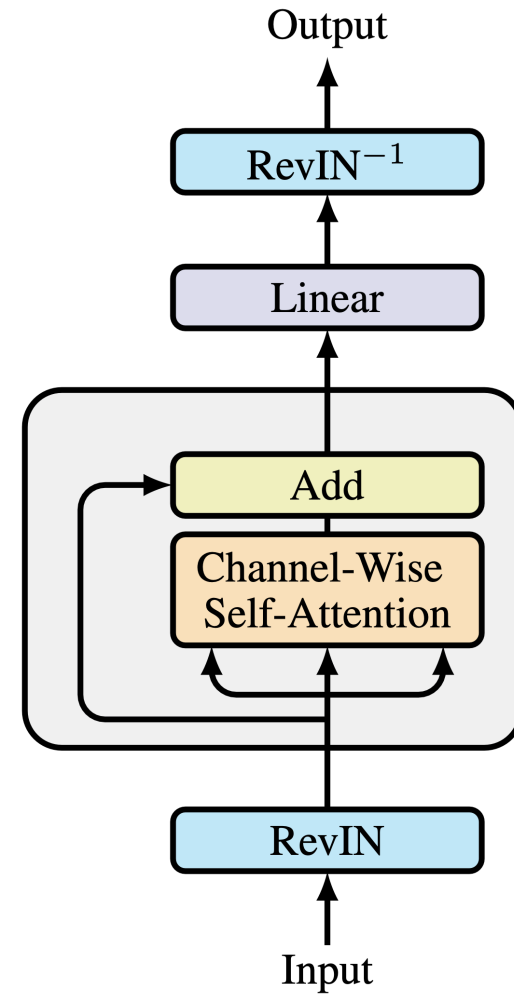
SAMformer (Ilbert et al. , ICML 2024)

A transformer-based Time Series Forecaster that actually works

SAMformer : Architecture

We can now design a simple model:

1. **RevIN layer** to be robust to train/test time shift
2. Shallow transformer with a **channel-wise attention**
3. We optimize it with **SAM**



SAMformer : Experimental results

SAMformer is better than all transformer-based models...

	iTrans*	PatchTST	FED*	TSMixer	Auto*	Pyra*	LogTrans	In*
Improvement	3.94%	11.13%	12.36%	14.33%	22.65%	61.88%	70.88%	72.20%

* formers

SAMformer : Experimental results

SAMformer is better than all transformer-based models...

	iTrans*	PatchTST	FED*	TSMixer	Auto*	Pyra*	LogTrans	In*
Improvement	3.94%	11.13%	12.36%	14.33%	22.65%	61.88%	70.88%	72.20%

* formers

... and the best mixer-based model ...

	TSMixer
Improvement	14.33%

SAMformer : Experimental results

SAMformer is better than all transformer-based models...

	iTrans*	PatchTST	FED*	TSMixer	Auto*	Pyra*	LogTrans	In*
Improvement	3.94%	11.13%	12.36%	14.33%	22.65%	61.88%	70.88%	72.20%

* formers

... and the best mixer-based model ...

	TSMixer
Improvement	14.33%

... while being much more simpler

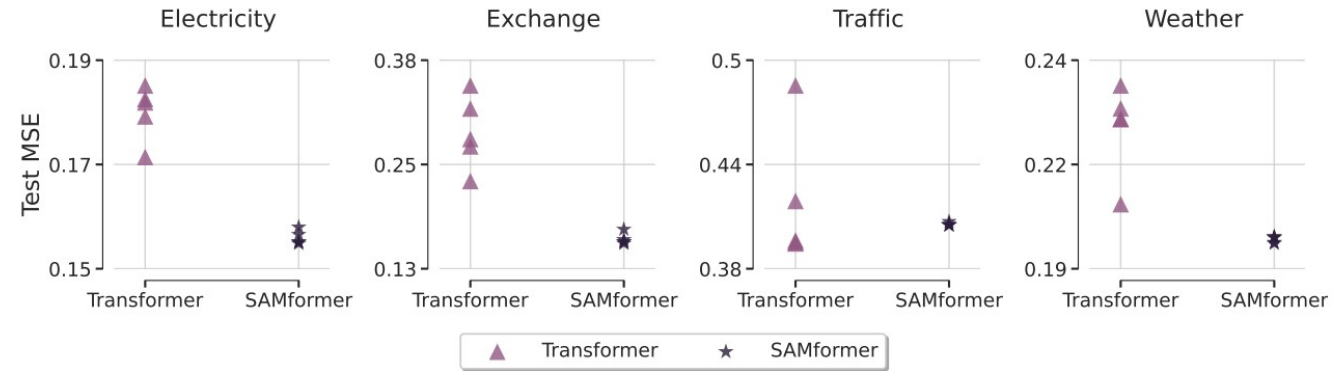
```
import torch
import torch.nn as nn
from models.utils import RevNorm

class SAMformer(nn.Module):
    def __init__(self, pred_len, num_heads=1, d_model=16):
        super(SAMformer, self).__init__()
        self.pred_len = pred_len
        self.rev_norm = RevNorm(axis=-2)
        self.attention_layer = nn.MultiheadAttention(embed_dim=d_model, num_heads=num_heads, batch_first=True)
        self.dense = nn.Linear(d_model, pred_len)

    def forward(self, x):
        x = self.rev_norm(x, mode='norm').transpose(1, 2) #RevIN normalization
        attention_output, _ = self.attention_layer(x, x, x) #Channel-Wise Attention Computation
        x = x + attention_output #Residual Connection
        x = self.dense(x).transpose(1, 2) #Final Linear Projection
        return self.rev_norm(x, mode='denorm') #RevIn denormalization
```

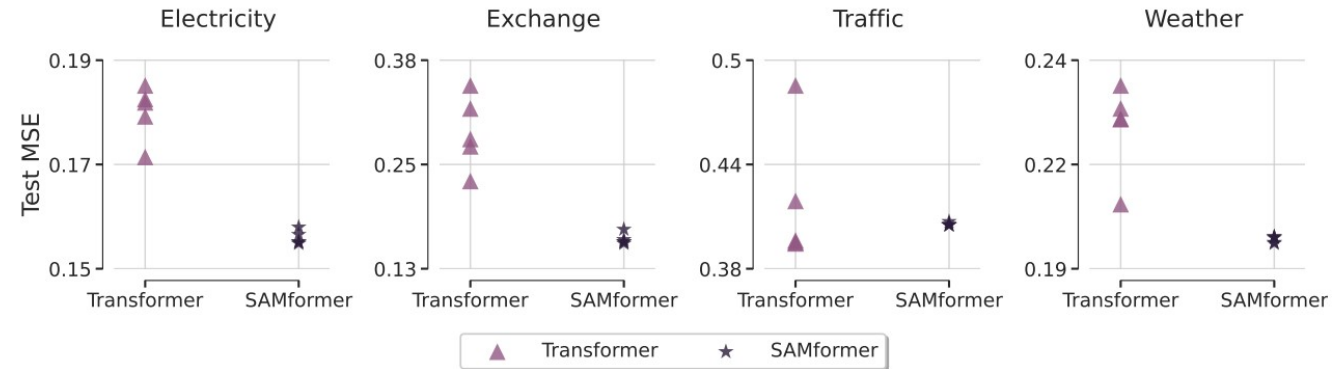
SAMformer : Experimental results

SAMformer is robust to random initialization ...



SAMformer : Experimental results

SAMformer is robust to random initialization ...

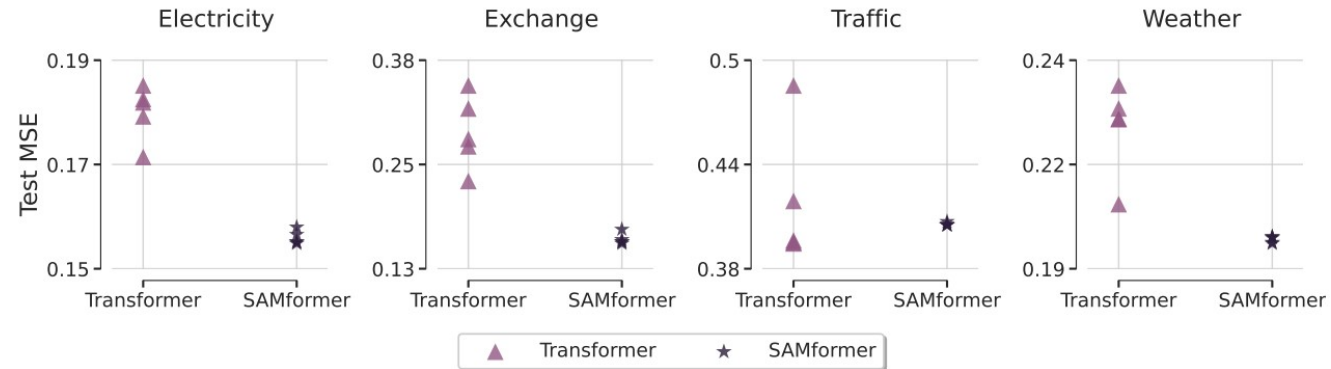


... and on par with the Foundation Model MOIRAI

Dataset	Full-shot	Zero-shot (Woo et al., 2024).		
	SAMformer	MOIRAI _{Small}	MOIRAI _{Base}	MOIRAI _{Large}
ETTh1	<u>0.410</u>	0.400	0.434	0.510
ETTh2	<u>0.344</u>	0.341	0.345	0.354
ETTm1	0.373	0.448	<u>0.381</u>	0.390
ETTm2	0.269	0.300	<u>0.272</u>	0.276
Electricity	0.181	0.233	<u>0.188</u>	<u>0.188</u>
Weather	0.260	<u>0.242</u>	0.238	0.259
Overall MSE improvement		6.9%	1.1%	7.6%

SAMformer : Experimental results

SAMformer is robust to random initialization ...



... and on par with the Foundation Model MOIRAI with many fewer parameters

Dataset	Full-shot	Zero-shot (Woo et al., 2024).		
	SAMformer	MOIRAI _{Small}	MOIRAI _{Base}	MOIRAI _{Large}
ETTh1	<u>0.410</u>	0.400	0.434	0.510
ETTh2	<u>0.344</u>	0.341	0.345	0.354
ETTm1	0.373	0.448	<u>0.381</u>	0.390
ETTm2	0.269	0.300	<u>0.272</u>	0.276
Electricity	0.181	0.233	<u>0.188</u>	<u>0.188</u>
Weather	0.260	<u>0.242</u>	0.238	0.259
Overall MSE improvement		6.9%	1.1%	7.6%

Model	# of parameters
MOIRAI _{small}	14M
MOIRAI _{base}	91M
MOIRAI _{large}	311M
SAMformer	280K

SAM vs σ – Reparametrization : Comparison

σ – **Reparametrization** (Zhen et al. 2023)

- smoothes the attention matrix with attention weights reparametrization

$$\hat{\mathbf{W}} = \frac{\gamma}{\|\mathbf{W}\|_2} \mathbf{W}$$

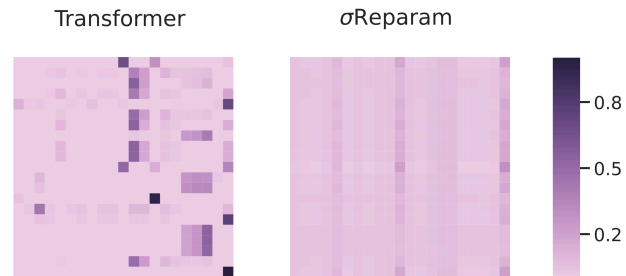
SAM vs σ – Reparametrization : Comparison

σ – Reparametrization (Zhen et al. 2023)

- smooths the attention matrix with attention weights reparametrization

$$\hat{\mathbf{W}} = \frac{\gamma}{\|\mathbf{W}\|_2} \mathbf{W}$$

- proved to be efficient in NLP ... but does not work for TS Forecasting



Observations:

- Transformers ignore diagonal elements
- σ – Reparametrization over-smoothes the attention matrix

} Rank Collapse
= Uninformative
Channel-Wise Attention

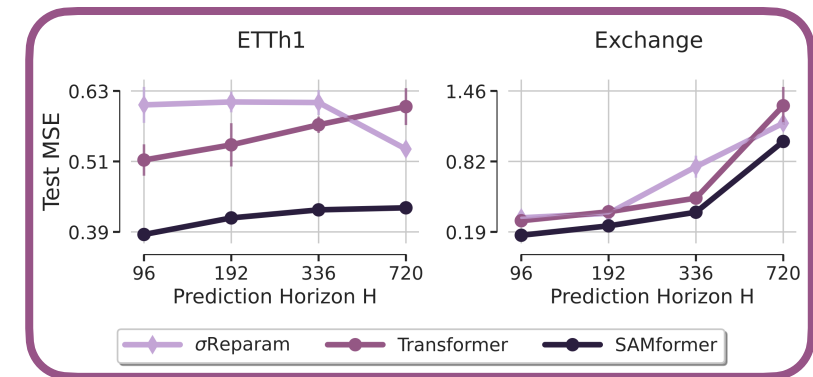
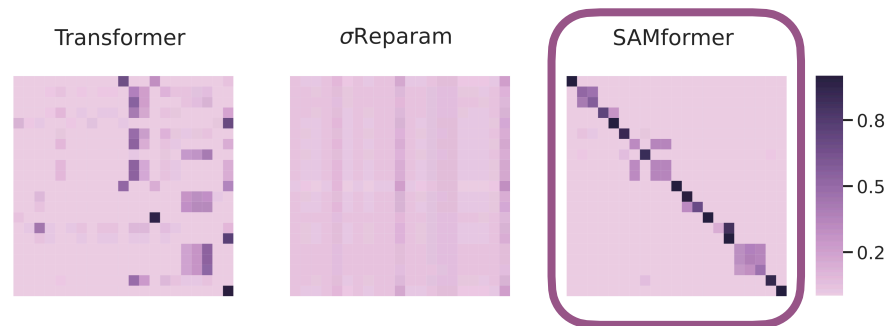
SAM vs σ – Reparametrization : Comparison

σ – Reparametrization (Zhen et al. 2023)

- smoothes the attention matrix with attention weights reparametrization

$$\hat{\mathbf{W}} = \frac{\gamma}{\|\mathbf{W}\|_2} \mathbf{W}$$

- proved to be efficient in NLP ... but does not work for TS Forecasting



Observations:

- Transformers ignore diagonal elements
- σ – Reparametrization over-smoothes the attention matrix
- SAMformer encourages feature self-correlation

Rank Collapse
= Uninformative Channel-Wise Attention

SAMformer : Ablation Studies

Ablation study on channel-wise attention and identity weight matrix attention ?

- Candidate 1: SAMformer with **temporal** attention (as used in all other transformers)

Model	Metrics	H	ETTh1	ETTh2	ETTm1	ETTm2	Electricity	Exchange	Traffic	Weather	Overall Improvement
Temporal Attention	MSE	96	0.496 \pm 0.009	0.401 \pm 0.011	0.542 \pm 0.063	0.330 \pm 0.034	0.291 \pm 0.025	0.684 \pm 0.218	0.933 \pm 0.188	0.225 \pm 0.005	12.97%
		192	0.510 \pm 0.014	0.414 \pm 0.020	0.615 \pm 0.056	0.394 \pm 0.033	0.294 \pm 0.024	0.434 \pm 0.063	0.647 \pm 0.131	0.254 \pm 0.001	
		336	0.549 \pm 0.017	0.396 \pm 0.014	0.620 \pm 0.046	0.436 \pm 0.081	0.290 \pm 0.016	0.473 \pm 0.014	0.656 \pm 0.113	0.292 \pm 0.000	
		720	0.604 \pm 0.017	0.396 \pm 0.010	0.694 \pm 0.055	0.469 \pm 0.005	0.307 \pm 0.014	1.097 \pm 0.084	-	0.346 \pm 0.000	
Temporal Attention	MAE	96	0.488 \pm 0.007	0.434 \pm 0.006	0.525 \pm 0.040	0.393 \pm 0.020	0.386 \pm 0.014	0.589 \pm 0.096	0.598 \pm 0.072	0.277 \pm 0.004	18.09%
		192	0.492 \pm 0.010	0.443 \pm 0.015	0.566 \pm 0.032	0.421 \pm 0.019	0.385 \pm 0.014	0.498 \pm 0.033	0.467 \pm 0.072	0.294 \pm 0.001	
		336	0.517 \pm 0.012	0.440 \pm 0.012	0.550 \pm 0.024	0.443 \pm 0.039	0.383 \pm 0.009	0.517 \pm 0.008	0.469 \pm 0.070	0.320 \pm 0.000	
		720	0.556 \pm 0.009	0.442 \pm 0.006	0.584 \pm 0.027	0.459 \pm 0.004	0.396 \pm 0.012	0.782 \pm 0.041	-	0.356 \pm 0.000	

SAMformer : Ablation Studies

Ablation study on channel-wise attention and identity weight matrix attention ?

- Candidate 1: SAMformer with **temporal** attention (as used in all other transformers)

Model	Metrics	H	ETTh1	ETTh2	ETTm1	ETTm2	Electricity	Exchange	Traffic	Weather	Overall Improvement
Temporal Attention	MSE	96	0.496 \pm 0.009	0.401 \pm 0.011	0.542 \pm 0.063	0.330 \pm 0.034	0.291 \pm 0.025	0.684 \pm 0.218	0.933 \pm 0.188	0.225 \pm 0.005	12.97%
		192	0.510 \pm 0.014	0.414 \pm 0.020	0.615 \pm 0.056	0.394 \pm 0.033	0.294 \pm 0.024	0.434 \pm 0.063	0.647 \pm 0.131	0.254 \pm 0.001	
		336	0.549 \pm 0.017	0.396 \pm 0.014	0.620 \pm 0.046	0.436 \pm 0.081	0.290 \pm 0.016	0.473 \pm 0.014	0.656 \pm 0.113	0.292 \pm 0.000	
		720	0.604 \pm 0.017	0.396 \pm 0.010	0.694 \pm 0.055	0.469 \pm 0.005	0.307 \pm 0.014	1.097 \pm 0.084	-	0.346 \pm 0.000	
	MAE	96	0.488 \pm 0.007	0.434 \pm 0.006	0.525 \pm 0.040	0.393 \pm 0.020	0.386 \pm 0.014	0.589 \pm 0.096	0.598 \pm 0.072	0.277 \pm 0.004	18.09%
		192	0.492 \pm 0.010	0.443 \pm 0.015	0.566 \pm 0.032	0.421 \pm 0.019	0.385 \pm 0.014	0.498 \pm 0.033	0.467 \pm 0.072	0.294 \pm 0.001	
		336	0.517 \pm 0.012	0.440 \pm 0.012	0.550 \pm 0.024	0.443 \pm 0.039	0.383 \pm 0.009	0.517 \pm 0.008	0.469 \pm 0.070	0.320 \pm 0.000	
		720	0.556 \pm 0.009	0.442 \pm 0.006	0.584 \pm 0.027	0.459 \pm 0.004	0.396 \pm 0.012	0.782 \pm 0.041	-	0.356 \pm 0.000	

- Candidate 2: SAMformer with **identity weight matrix** attention

Model	Metrics	H	ETTh1	ETTh2	ETTm1	ETTm2	Electricity	Exchange	Traffic	Weather	Overall Improvement
Identity Attention	MSE	96	0.477 \pm 0.059	0.346 \pm 0.055	0.345 \pm 0.027	0.201 \pm 0.035	0.175 \pm 0.015	0.179 \pm 0.031	0.416 \pm 0.037	0.206 \pm 0.019	11.93%
		192	0.467 \pm 0.074	0.374 \pm 0.031	0.384 \pm 0.042	0.248 \pm 0.016	0.189 \pm 0.022	0.320 \pm 0.070	0.437 \pm 0.041	0.236 \pm 0.002	
		336	0.512 \pm 0.070	0.372 \pm 0.024	0.408 \pm 0.032	0.303 \pm 0.022	0.211 \pm 0.019	0.443 \pm 0.071	0.500 \pm 0.155	0.277 \pm 0.003	
		720	0.505 \pm 0.107	0.405 \pm 0.012	0.466 \pm 0.043	0.397 \pm 0.029	0.233 \pm 0.019	1.123 \pm 0.076	0.468 \pm 0.021	0.338 \pm 0.009	
	MAE	96	0.473 \pm 0.041	0.395 \pm 0.033	0.376 \pm 0.019	0.294 \pm 0.027	0.283 \pm 0.023	0.320 \pm 0.023	0.301 \pm 0.039	0.259 \pm 0.021	4.18%
		192	0.463 \pm 0.055	0.413 \pm 0.022	0.399 \pm 0.030	0.321 \pm 0.012	0.291 \pm 0.029	0.418 \pm 0.043	0.314 \pm 0.042	0.278 \pm 0.002	
		336	0.490 \pm 0.049	0.413 \pm 0.015	0.411 \pm 0.019	0.354 \pm 0.018	0.309 \pm 0.021	0.498 \pm 0.041	0.350 \pm 0.106	0.305 \pm 0.003	
		720	0.496 \pm 0.066	0.438 \pm 0.008	0.444 \pm 0.030	0.406 \pm 0.017	0.322 \pm 0.021	0.788 \pm 0.021	0.325 \pm 0.023	0.347 \pm 0.009	

Conclusions on SAMformer (ILBERT et al, ICML oral 2024)

1. We studied **pitfalls of transformers** in time series forecasting
 - Sharp loss landscape = lack of generalization
2. Our proposal **SAMformer**
 - **SAMformer** = RevIN + channel-wise attention + SAM optimization
 - **SOTA** in long-term multivariate time series forecasting
 - **Consistent** = same architecture of different horizons/datasets
 - **Lightweight** = the smallest SOTA model
 - **On par with** the large foundation model **MOIRAI**
3. Can inspire further work to enhance our simple architecture.

Thank you

Check my website for code, paper and slides



Let's discuss potential post-PhD opportunities