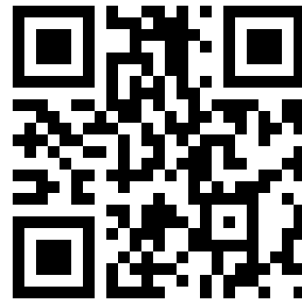# SAMformer : Unlocking the potential of Transformers in Time-Series Forecasting

## Presentation : Romain ILBERT*
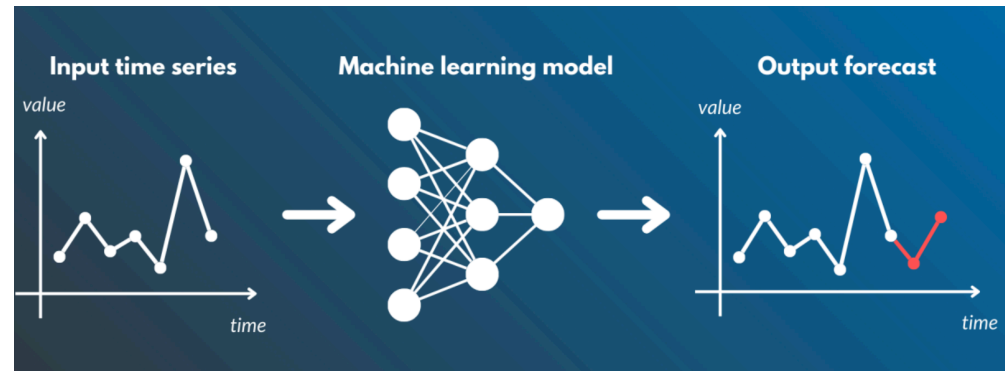
Check my website for code and paper :

* PhD Student at LIPADE, Paris Descartes University & Huawei Paris Research Center

Accepted as an Oral at ICML 2024

# Time Series Forecasting : A Definition

**Problem Setup**
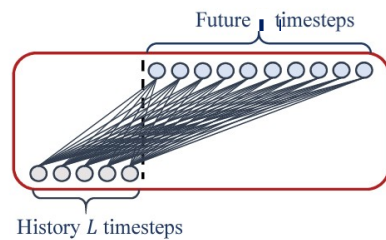
1. Given past observations, predict future ones



2. Univariate (single channel) vs. multivariate (multi-channels)

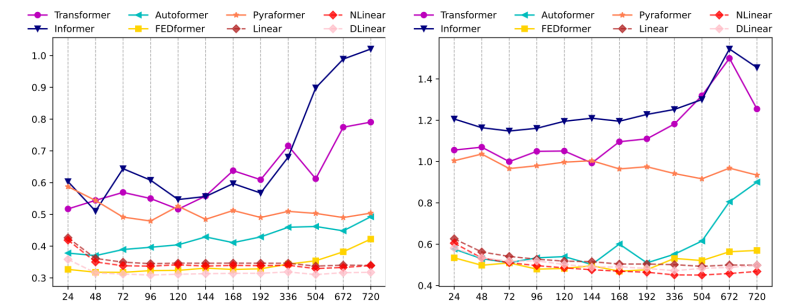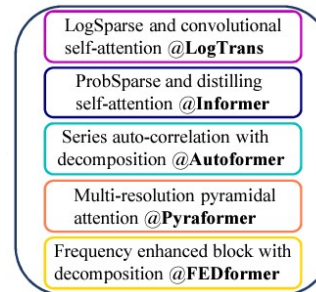3. Short, medium and long-term horizon

# Failure of Transformers

**Motivation :**

1. Consider the simplest linear model for forecasting
2. Compare it to state-of-the-art transformers



(a) **24** steps-**ETTh1**      (b) **720** steps-**ETTh1**

This Linear model surpasses the SOTA FEDformer (ICML'22) in most cases by 20%~50%

**Main conclusions from *Zeng et al.***

1. Transformer-based methods don't work well in forecasting
2. Embarrassing failure in most basic scenario

… yet they <u>dominate NLP and vision</u>. Why?

1. Are Transformers Effective for Time Series Forecasting ? Zeng et al. 2023.

# SAMformer (Ilbert et al. , ICML Oral 2024)

A transformer-based TS forecaster that actually works
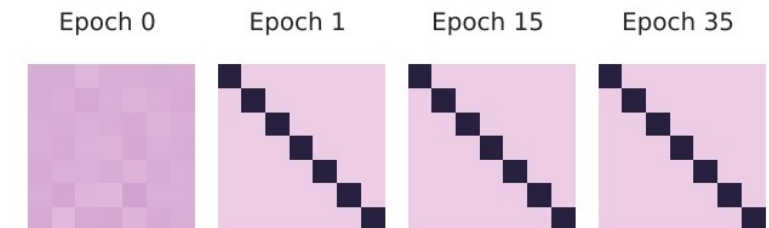
# Simple Toy Regression Example

**Context :**

- Consider a simple linear regression problem

- Transformer defined as $f(\mathbf{X}) = [\mathbf{X} + A(\mathbf{X})\mathbf{X}\mathbf{W}_V\mathbf{W}_O]\mathbf{W}$

  with **channel-wise attention** (DxD matrix, rather than LxL)

**Conclusions :**

1. Linear Transformer severely overfits

2. And works better if we freeze the attention

3. Because the attention get stuck at the identity

   matrix and does not move afterwards

**Pathological behavior suggesting sharp local minima**

$$\mathbf{Y} = \mathbf{X}\mathbf{W}_{\text{toy}} + \epsilon \quad \text{(L=512, H=96, D=7)}$$





Epoch 0    Epoch 1    Epoch 15    Epoch 35

# Simple Toy Regression Example

**Why transformers fail?**

1. Transformers have a <u>sharp loss landscape</u> and suffer from <u>entropy collapse</u>



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

3. And no changing in the optimizer helps to solve this



(a) SGD and AdamW with wd = 1e−3          (b) AdamW with wd ∈ {1e−5, 1e−4}.

# Simple Toy Regression Example : To fix the sharpness issue

1. $\sigma -$ **Reparametrization** (Zhen et al. 2023)

   - Make attention matrix more uniform to avoid entropy collapse

$$\hat{\mathbf{W}} = \frac{\gamma}{\|\mathbf{W}\|_2}\mathbf{W}$$
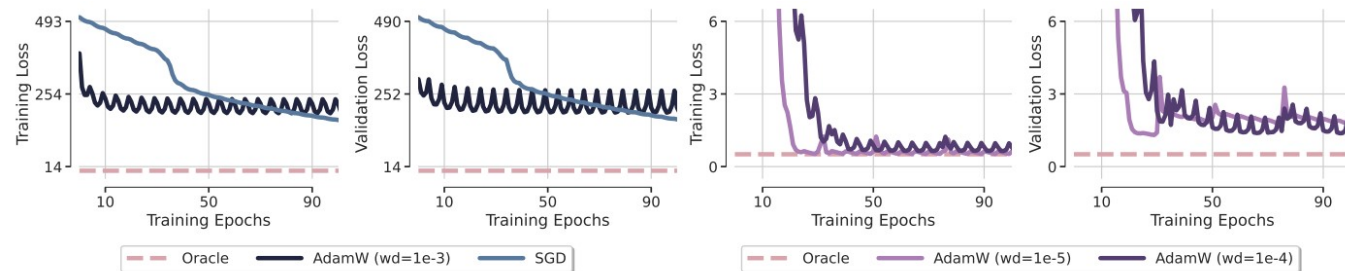
2. **Sharpness - Aware Minimization** (Foret et al. 2021, Chen et al. 2022)

   - Converge toward weights that lie in neighborhoods having uniformly low loss

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

**SAM = desired solution**

# Congrats you now know how to solve

# a linear regression problem with transformers!

# SAMformer Architecture (ILBERT et al, ICML oral 2024)

**Let's put it all together now:**

1. Shallow transformer with a **channel-wise attention**

2. **RevIN layer** to be robust to train/test time shift

3. We optimize it with **SAM**

**SAMformer = 15 lines of code**

# Experimental results (ILBERT et al, ICML oral 2024)

1. Datasets

| Dataset | ETTh1/ETTh2 | ETTm1/ETTm2 | Electricity | Exchange | Traffic | Weather |
|---|---|---|---|---|---|---|
| # features | 7 | 7 | 321 | 8 | 862 | 21 |
| # time steps | 17420 | 69680 | 26304 | 7588 | 17544 | 52696 |
| Granularity | 1 hour | 15 minutes | 1 hour | 1 day | 1 hour | 10 minutes |

2. Baselines

  - TSmixer: MLPmixer model from Google (SOTA in 2023)

  - Transformers: iTransformer (ICLR'24), PatchTST (ICLR'23), FEDformer (ICML'22), Pyraformer (ICLR'22), Informer (AAAI'21), Autoformer (NeurIPS'21), LogTrans (NeurIPS'19)

3. Nbr of parameters

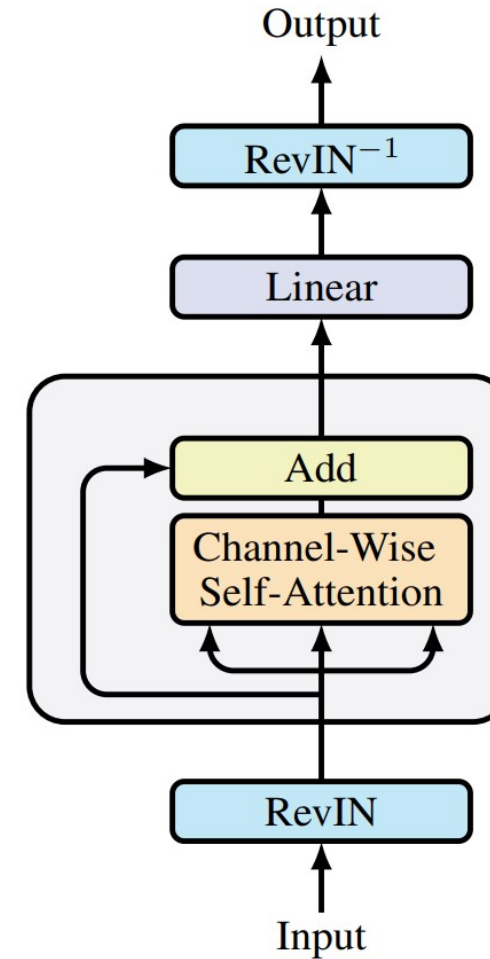  - SAMFormer is **smaller** and **more consistent** than TSMixer. The **same model** for all datasets/horizons

  - Avg Ratio = nbre params TSMixer / nbre params SAMFormer

| Dataset | $H = 96$ | | $H = 192$ | | $H = 336$ | | $H = 720$ | | Total |
|---|---|---|---|---|---|---|---|---|---|
| | SAMformer | TSMixer | SAMformer | TSMixer | SAMformer | TSMixer | SAMformer | TSMixer | |
| ETT | 50272 | 124142 | 99520 | 173390 | 173392 | 247262 | 369904 | 444254 | - |
| Exchange | 50272 | 349344 | 99520 | 398592 | 173392 | 472464 | 369904 | 669456 | - |
| Weather | 50272 | 121908 | 99520 | 171156 | 173392 | 245028 | 369904 | 442020 | - |
| Electricity | 50272 | 280676 | 99520 | 329924 | 173392 | 403796 | 369904 | 600788 | - |
| Traffic | 50272 | 793424 | 99520 | 842672 | 173392 | 916544 | 369904 | 1113536 | - |
| **Avg. Ratio** | 6.64 | | 3.85 | | 2.64 | | 1.77 | | **3.73** |

# Experimental results (ILBERT et al, ICML oral 2024)

| | with SAM | | without SAM | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | **SAMformer** | TSMixer | Transformer | TSMixer | In* | Auto* | FED* | Pyra[†] | LogTrans[†] |
| **Overall MSE improvement** | 5.25% | 16.96% | 14.33% | 72.20% | 22.65% | 12.36% | 61.88% | 70.88% | |

1. SAMformer is **14% better** than TSMixer, **11.13% better** than PatchTST, **3.94%** better than iTransformer
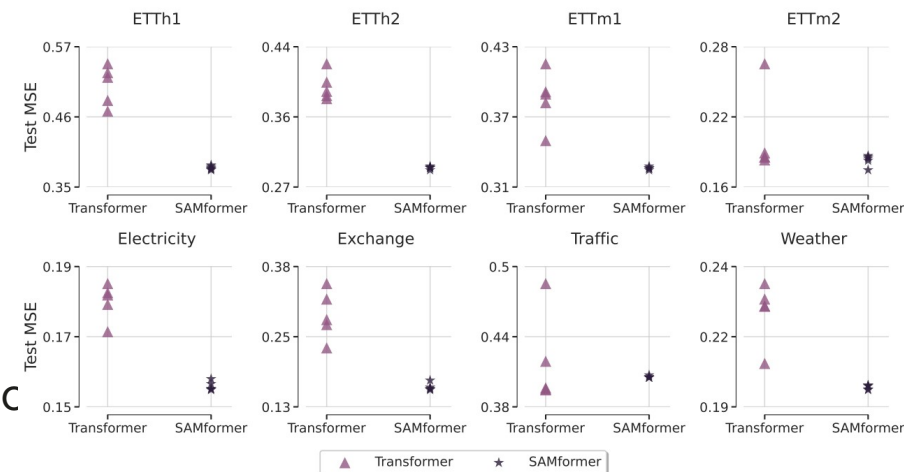
   - much better than all transformer-based models

2. <u>Sharpness-aware minimization</u> **improves** TSMixer as well

3. SAMformer is robust to random initialization

   - Very low variance for different random seeds compared to

Transformer

# Experimental results (ILBERT et al, ICML oral 2024)

SAMFormer is **on par with MOIRAI foundation model**

- MORAI (Salesforce + Singapore University)

- trained on LOTSA with <u>27B samples</u> from <u>9 domains</u>

- comes in 3 sizes: small (14M), base (91M) and Large (311M)

| | | MOIRAI$_{Small}$ | MOIRAI$_{Base}$ | MOIRAI$_{Large}$ | SAMformer |
|---|---|---|---|---|---|
| ETTh1 | MSE | **0.400** | <u>0.434</u> | 0.510 | 0.41 |
| | MAE | **0.424** | <u>0.438</u> | 0.469 | |
| ETTh2 | MSE | **0.341** | <u>0.345</u> | 0.354 | 0.344 |
| | MAE | <u>0.379</u> | 0.382 | **0.376** | |
| ETTm1 | MSE | 0.448 | **0.381** | 0.390 | **0.373** |
| | MAE | 0.409 | **0.388** | <u>0.389</u> | |
| ETTm2 | MSE | 0.300 | **0.272** | <u>0.276</u> | **0.2685** |
| | MAE | 0.341 | <u>0.321</u> | **0.320** | |
| Electricity | MSE | 0.233 | 0.188 | <u>0.188</u> | **0.181** |
| | MAE | 0.320 | 0.274 | <u>0.273</u> | |
| Weather | MSE | <u>0.242</u> | **0.238** | 0.259 | 0.26 |
| | MAE | <u>0.267</u> | **0.261** | 0.275 | |

# SAM vs $\sigma -$ Reparametrization

## Comparison with $\sigma -$ **Reparametrization**

- proved to be efficient in NLP ... but didn't work for us



(a) Comparison of `Transformer`, $\sigma$Reparam and `SAMformer`. (b) Comparison of `SAMformer` and `SAMformer` + $\sigma$Reparam.

## Observations:

- Transformers ignores diagonal elements

- SAMformer strongly encourages feature self-correlation (as in ViTs)

- Weight reparametrization oversmoothes the attention matrix

# SAM vs $\sigma$ − Reparametrization

Oversmoothing = rank collapse

- we prove that

**Proposition 2.2** (Upper bound on the nuclear norm)
Let $\mathbf{X} \in \mathbb{R}^{D \times L}$ be an input sequence. Assuming $\mathbf{W}_Q \mathbf{W}_K^\top = \mathbf{W}_K \mathbf{W}_Q^\top \succcurlyeq 0$, we have

$$\|\mathbf{X}\mathbf{W}_Q\mathbf{W}_K^\top\mathbf{X}^\top\|_* \leq \|\mathbf{W}_Q\mathbf{W}_K^\top\|_2 \|\mathbf{X}\|_F^2.$$

Roughly = **rank** of the **attention matrix**

Minimized by reparametrization

- maximizing the entropy of the attention = rank collapse

- rank collapse = uninformative channel-wise attention

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

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

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

- Overall Improvement : Improvement of SAMFormer over both candidates

| Model | Metrics | H | ETTh1 | ETTh2 | ETTm1 | ETTm2 | Electricity | Exchange | Traffic | Weather | Overall Improvement |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Temporal Attention | MSE | 96 | $0.496_{\pm0.009}$ | $0.401_{\pm0.011}$ | $0.542_{\pm0.063}$ | $0.330_{\pm0.034}$ | $0.291_{\pm0.025}$ | $0.684_{\pm0.218}$ | $0.933_{\pm0.188}$ | $0.225_{\pm0.005}$ | |
| | | 192 | $0.510_{\pm0.014}$ | $0.414_{\pm0.020}$ | $0.615_{\pm0.056}$ | $0.394_{\pm0.033}$ | $0.294_{\pm0.024}$ | $0.434_{\pm0.063}$ | $0.647_{\pm0.131}$ | $0.254_{\pm0.001}$ | 12.97% |
| | | 336 | $0.549_{\pm0.017}$ | $0.396_{\pm0.014}$ | $0.620_{\pm0.046}$ | $0.436_{\pm0.081}$ | $0.290_{\pm0.016}$ | $0.473_{\pm0.014}$ | $0.656_{\pm0.113}$ | $0.292_{\pm0.000}$ | |
| | | 720 | $0.604_{\pm0.017}$ | $0.396_{\pm0.010}$ | $0.694_{\pm0.055}$ | $0.469_{\pm0.005}$ | $0.307_{\pm0.014}$ | $1.097_{\pm0.084}$ | - | $0.346_{\pm0.000}$ | |
| | MAE | 96 | $0.488_{\pm0.007}$ | $0.434_{\pm0.006}$ | $0.525_{\pm0.040}$ | $0.393_{\pm0.020}$ | $0.386_{\pm0.014}$ | $0.589_{\pm0.096}$ | $0.598_{\pm0.072}$ | $0.277_{\pm0.004}$ | |
| | | 192 | $0.492_{\pm0.010}$ | $0.443_{\pm0.015}$ | $0.566_{\pm0.032}$ | $0.421_{\pm0.019}$ | $0.385_{\pm0.014}$ | $0.498_{\pm0.033}$ | $0.467_{\pm0.072}$ | $0.294_{\pm0.001}$ | 18.09% |
| | | 336 | $0.517_{\pm0.012}$ | $0.440_{\pm0.012}$ | $0.550_{\pm0.024}$ | $0.443_{\pm0.039}$ | $0.383_{\pm0.009}$ | $0.517_{\pm0.008}$ | $0.469_{\pm0.070}$ | $0.320_{\pm0.000}$ | |
| | | 720 | $0.556_{\pm0.009}$ | $0.442_{\pm0.006}$ | $0.584_{\pm0.027}$ | $0.459_{\pm0.004}$ | $0.396_{\pm0.012}$ | $0.782_{\pm0.041}$ | - | $0.356_{\pm0.000}$ | |

| Model | Metrics | H | ETTh1 | ETTh2 | ETTm1 | ETTm2 | Electricity | Exchange | Traffic | Weather | Overall Improvement |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Identity Attention | MSE | 96 | $0.477_{\pm0.059}$ | $0.346_{\pm0.055}$ | $0.345_{\pm0.027}$ | $0.201_{\pm0.035}$ | $0.175_{\pm0.015}$ | $0.179_{\pm0.031}$ | $0.416_{\pm0.037}$ | $0.206_{\pm0.019}$ | |
| | | 192 | $0.467_{\pm0.074}$ | $0.374_{\pm0.031}$ | $0.384_{\pm0.042}$ | $0.248_{\pm0.016}$ | $0.189_{\pm0.022}$ | $0.320_{\pm0.070}$ | $0.437_{\pm0.041}$ | $0.236_{\pm0.002}$ | 11.93% |
| | | 336 | $0.512_{\pm0.070}$ | $0.372_{\pm0.024}$ | $0.408_{\pm0.032}$ | $0.303_{\pm0.022}$ | $0.211_{\pm0.019}$ | $0.443_{\pm0.071}$ | $0.500_{\pm0.155}$ | $0.277_{\pm0.003}$ | |
| | | 720 | $0.505_{\pm0.107}$ | $0.405_{\pm0.012}$ | $0.466_{\pm0.043}$ | $0.397_{\pm0.029}$ | $0.233_{\pm0.019}$ | $1.123_{\pm0.076}$ | $0.468_{\pm0.021}$ | $0.338_{\pm0.009}$ | |
| | MAE | 96 | $0.473_{\pm0.041}$ | $0.395_{\pm0.033}$ | $0.376_{\pm0.019}$ | $0.294_{\pm0.027}$ | $0.283_{\pm0.023}$ | $0.320_{\pm0.023}$ | $0.301_{\pm0.039}$ | $0.259_{\pm0.021}$ | |
| | | 192 | $0.463_{\pm0.055}$ | $0.413_{\pm0.022}$ | $0.399_{\pm0.030}$ | $0.321_{\pm0.012}$ | $0.291_{\pm0.029}$ | $0.418_{\pm0.043}$ | $0.314_{\pm0.042}$ | $0.278_{\pm0.002}$ | 4.18% |
| | | 336 | $0.490_{\pm0.049}$ | $0.413_{\pm0.015}$ | $0.411_{\pm0.019}$ | $0.354_{\pm0.018}$ | $0.309_{\pm0.021}$ | $0.498_{\pm0.041}$ | $0.350_{\pm0.106}$ | $0.305_{\pm0.003}$ | |
| | | 720 | $0.496_{\pm0.066}$ | $0.438_{\pm0.008}$ | $0.444_{\pm0.030}$ | $0.406_{\pm0.017}$ | $0.322_{\pm0.021}$ | $0.788_{\pm0.021}$ | $0.325_{\pm0.023}$ | $0.347_{\pm0.009}$ | |

# 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 large foundation model MORAI

Thank you.