



ICLR



SNU IMLAB

Dept. of Industrial Engineering
Seoul National University

MambaSL : Exploring Single-Layer Mamba for Time Series Classification

Yoo-Min Jung,

pamela7384@gmail.com

Leekyung Kim

klk97@snu.ac.kr

Department of Industrial Engineering, Seoul National University

Motivation

► Motivation

1. Underexplored in time series classification (TSC) using Mamba

- in time series forecasting (TSF)
 - TimeMachine (Ahamed & Cheng, 2024), S-Mamba (Wang et al., 2025), ...
 - showed promising performance with efficiency
- in TSC
 - poor performance with vanilla Mamba (Wang et al., 2024)
 - TSCMamba (Ahamed & Cheng, 2025)
 - only variant that is tested on extensive benchmark datasets

→ Find out **why** naively applying Mamba to TSC is not working

► Motivation

2. Issues in benchmarking TSC results

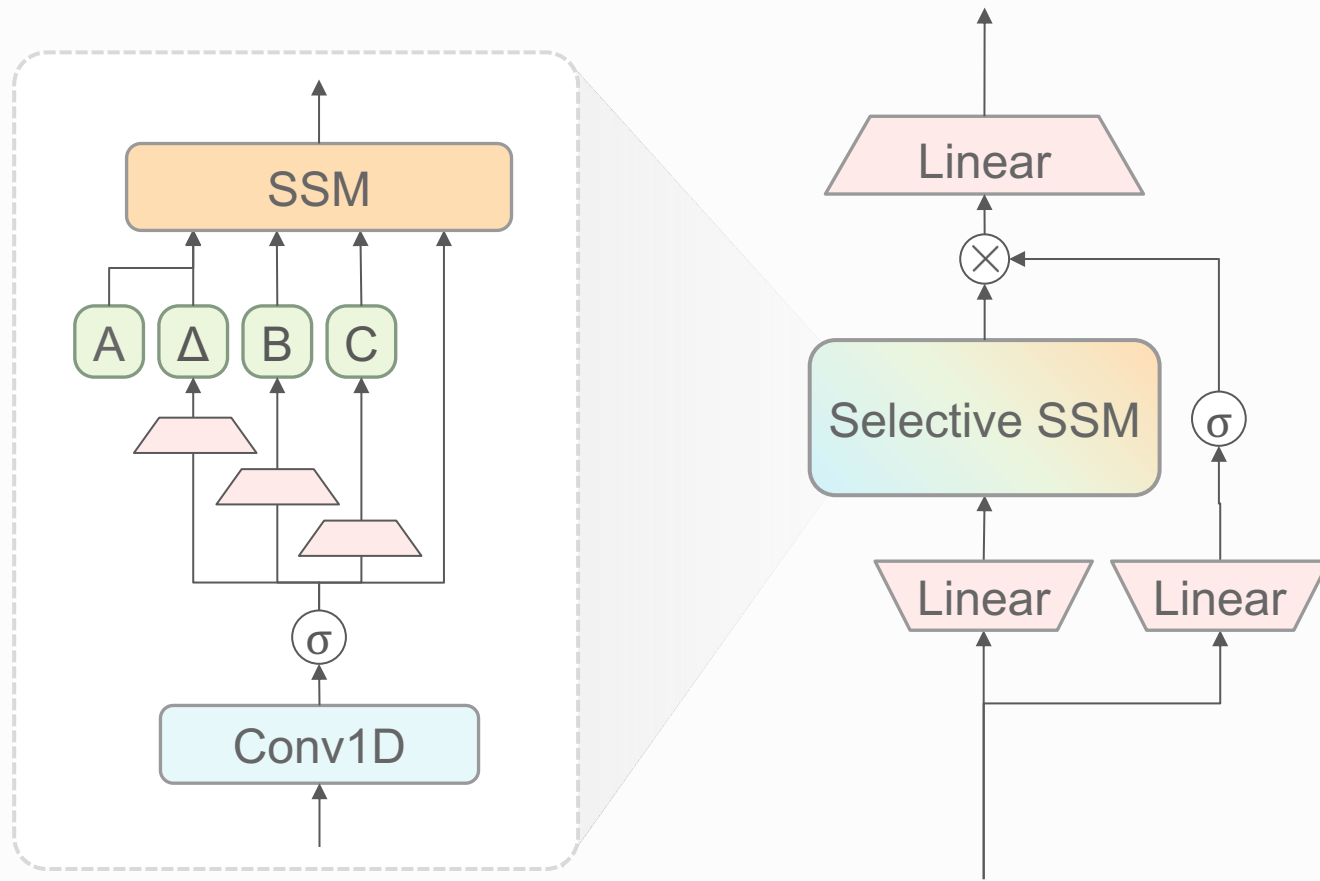
- coverage
 - each study uses a different subset of the UEA dataset (Bagnall et al., 2018)
- fairness
 - uneven hyperparameter tuning (especially for TSF-origin model)
- reproducibility
 - reported results often fail to reproduce original results (Eldele et al., 2024)

→ establish fair, reproducible, and full UEA 30 benchmarking

Background

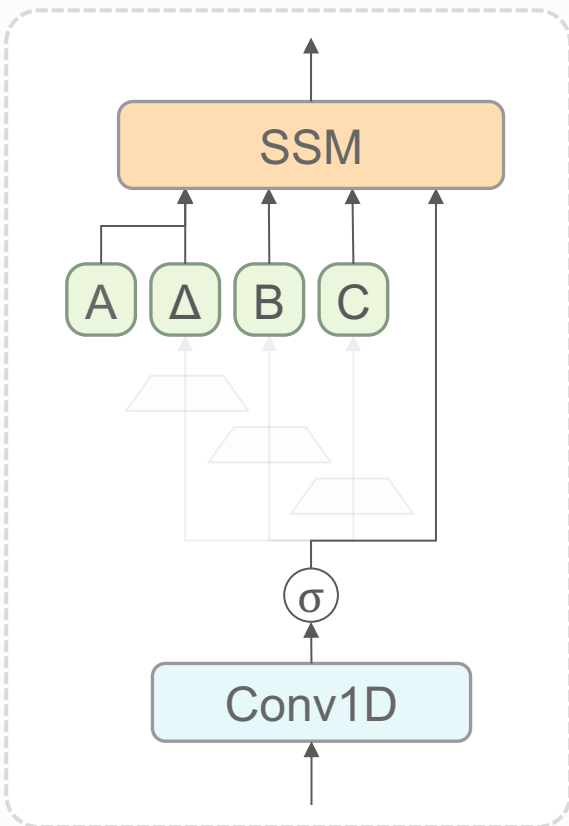
► Background

- **Mamba** = Selective SSM + Gated Linear Unit



Background

- SSM (state space model)



✓ $s(t)$: state vector, $u(t), y(t)$: input/output scalar

$$\dot{s}(t) = \mathbf{A}s(t) + \mathbf{B}u(t), \quad y(t) = \mathbf{C}s(t) + \mathbf{D}u(t), \quad (5)$$

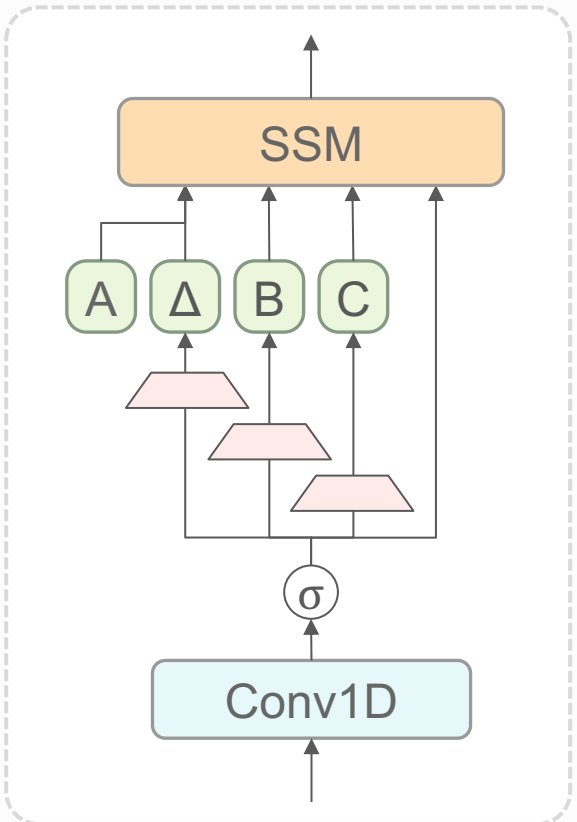
$$s_t = \bar{\mathbf{A}}s_{t-1} + \bar{\mathbf{B}}u_t, \quad y_t = \mathbf{C}s_t + \mathbf{D}u_t, \quad (6)$$

Discretize
(fixed step $\Delta > 0$)

- Characteristics
 - LTI (Linear Time-Invariant) System
 - Channel Independent

Background

- Selective SSM



✓ $\tilde{\mathbf{x}}_t$: given input at time $t \rightarrow \Delta_t, \mathbf{B}_t, \mathbf{C}_t$

$$\begin{aligned} \Delta_t &= \phi_{\Delta}(\tilde{\mathbf{x}}_t) = \zeta(\text{Linear}_{d_m}^{\text{bias}}(\text{Linear}_{d_r}(\tilde{\mathbf{x}}_t))) \in \mathbb{R}_{>0}^{d_m}, \\ \mathbf{B}_t &= \phi_B(\tilde{\mathbf{x}}_t) = \text{Linear}_{d_s}(\tilde{\mathbf{x}}_t) \in \mathbb{R}^{d_s \times 1}, \\ \mathbf{C}_t &= \phi_C(\tilde{\mathbf{x}}_t) = \text{Linear}_{d_s}(\tilde{\mathbf{x}}_t)^{\top} \in \mathbb{R}^{1 \times d_s}. \end{aligned} \quad (7)$$

- Characteristics
 - Time-Varying (TV) System
 - Channel Mixed

Background

- On time variance in Δ , B , and C

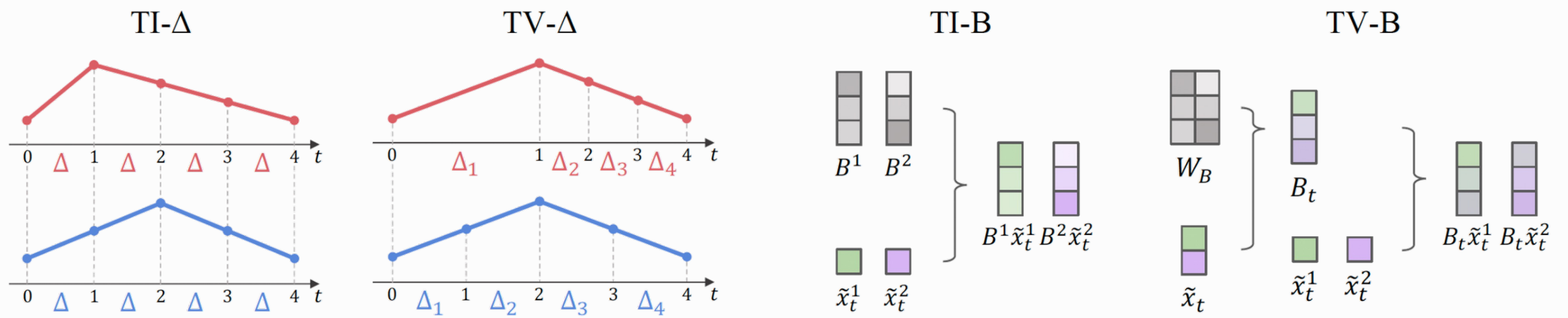
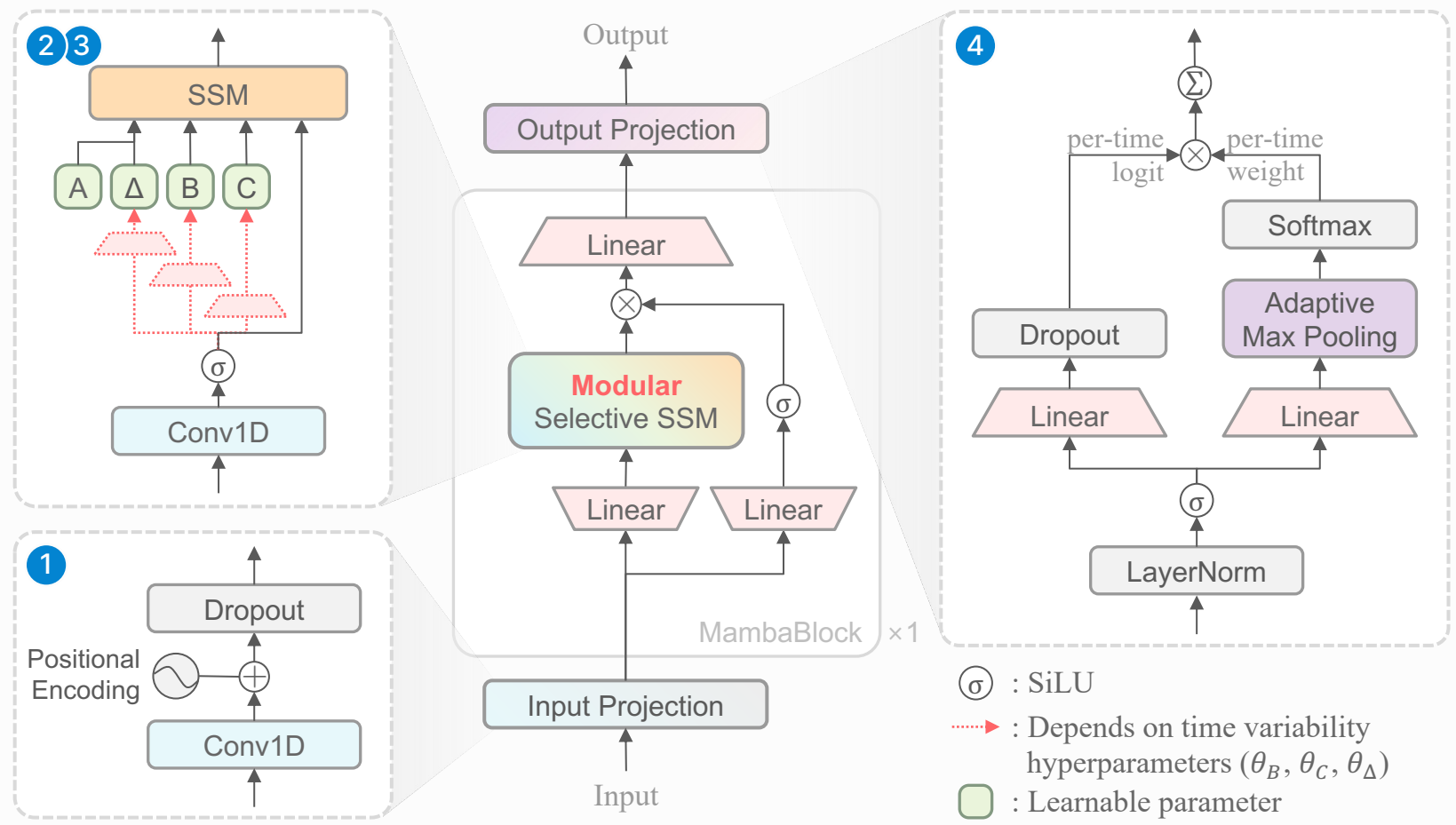


Figure 1: TI/TV parameterization of (left) Δ and (right) B in the SSM. TI- Δ fixes the update rate, while TV- Δ adapts it to align sequences with varying speeds. TI- B preserves channel independence, whereas TV- B introduces input-dependent mixing. C follows the same pattern as B at the output stage, highlighting the temporal pacing of Δ versus the spatial routing of B and C .

Method & Results

Method

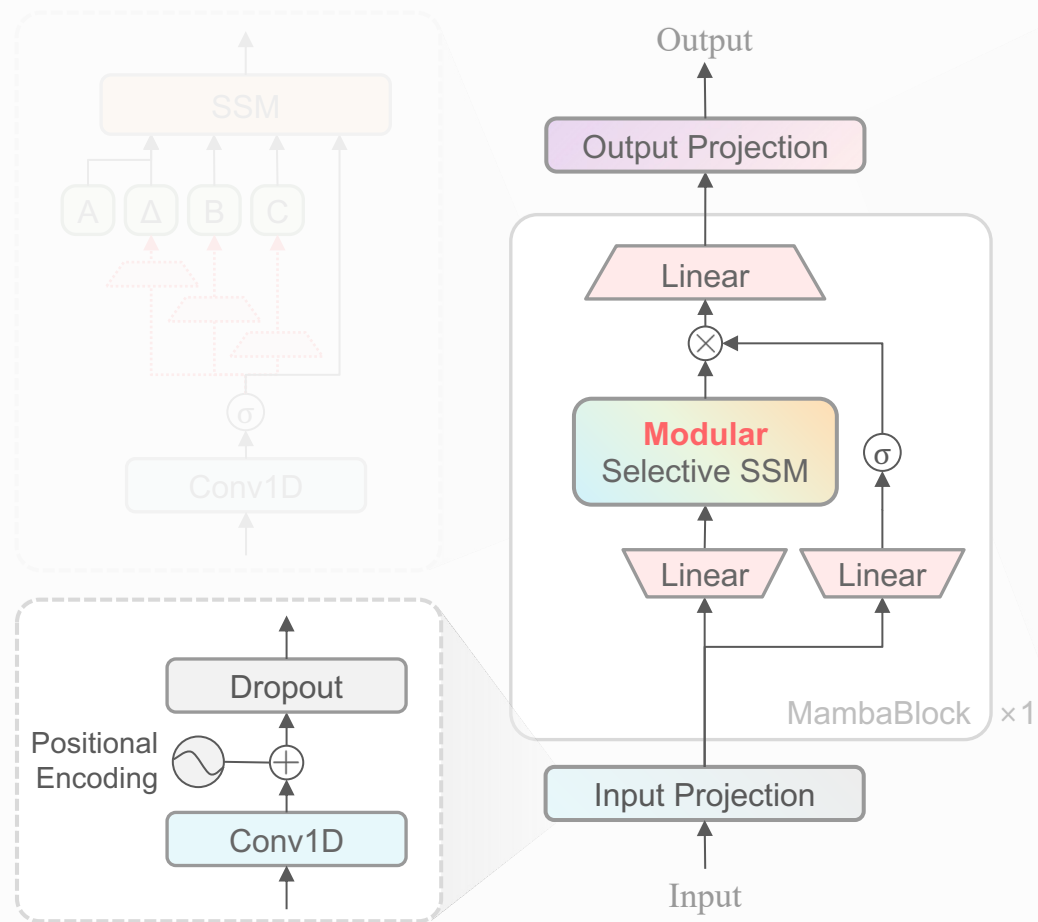
- **MambaSL** : TSC framework with a single-layer Mamba



- 1 Scale input projection
- 2 Modularize time (in)variance
- 3 Remove skip connection
- 4 Aggregate via adaptive pooling

Method

- (H1) Scale Input Projection



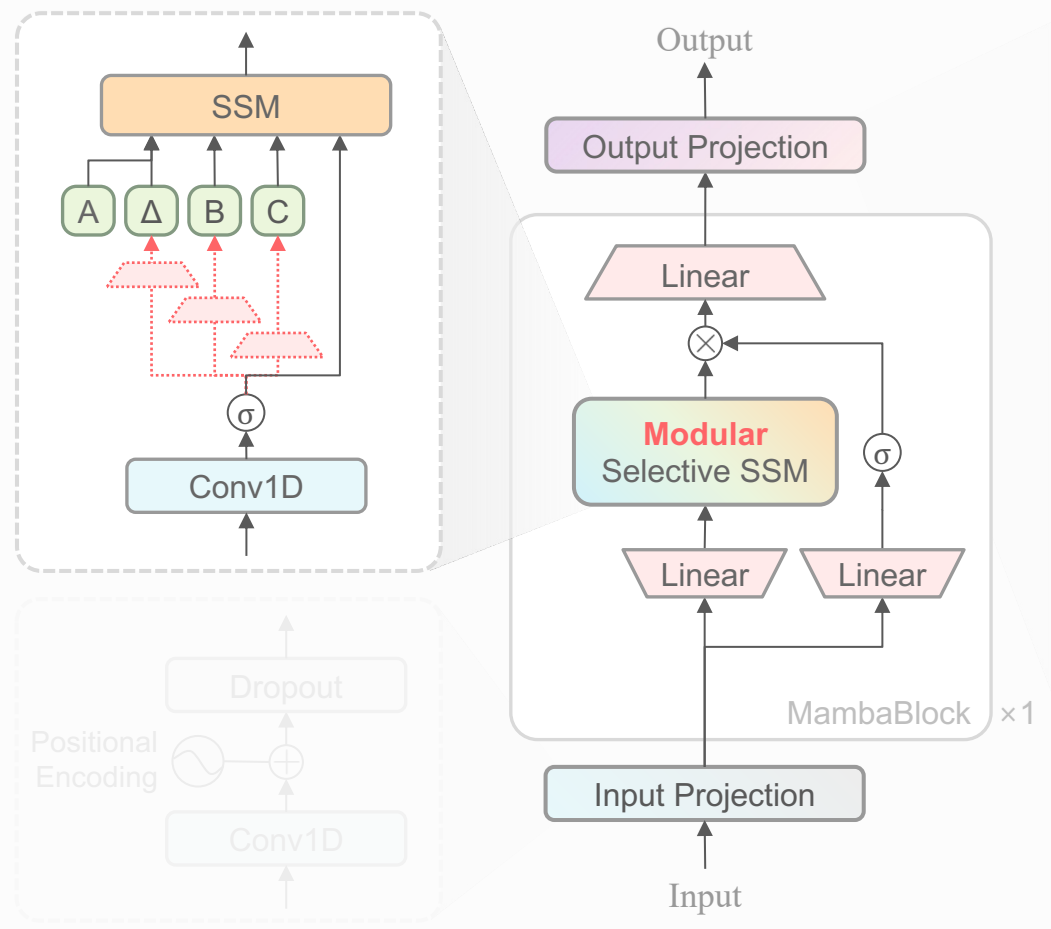
$$k = \max(k_{\min}, \lfloor \lambda L \rfloor) \quad (10)$$

- Mamba's gating mechanism modulates the SSM output based on input projection.
- time series dataset may require larger receptive fields to generate token embedding
 - especially for dataset with high sampling rate
 - especially for dataset with long sequence length

σ : SiLU
→ : Depends on time variability hyperparameters ($\theta_B, \theta_C, \theta_\Delta$)
 : Learnable parameter

Method

- (H2) Modularize time (in)variance

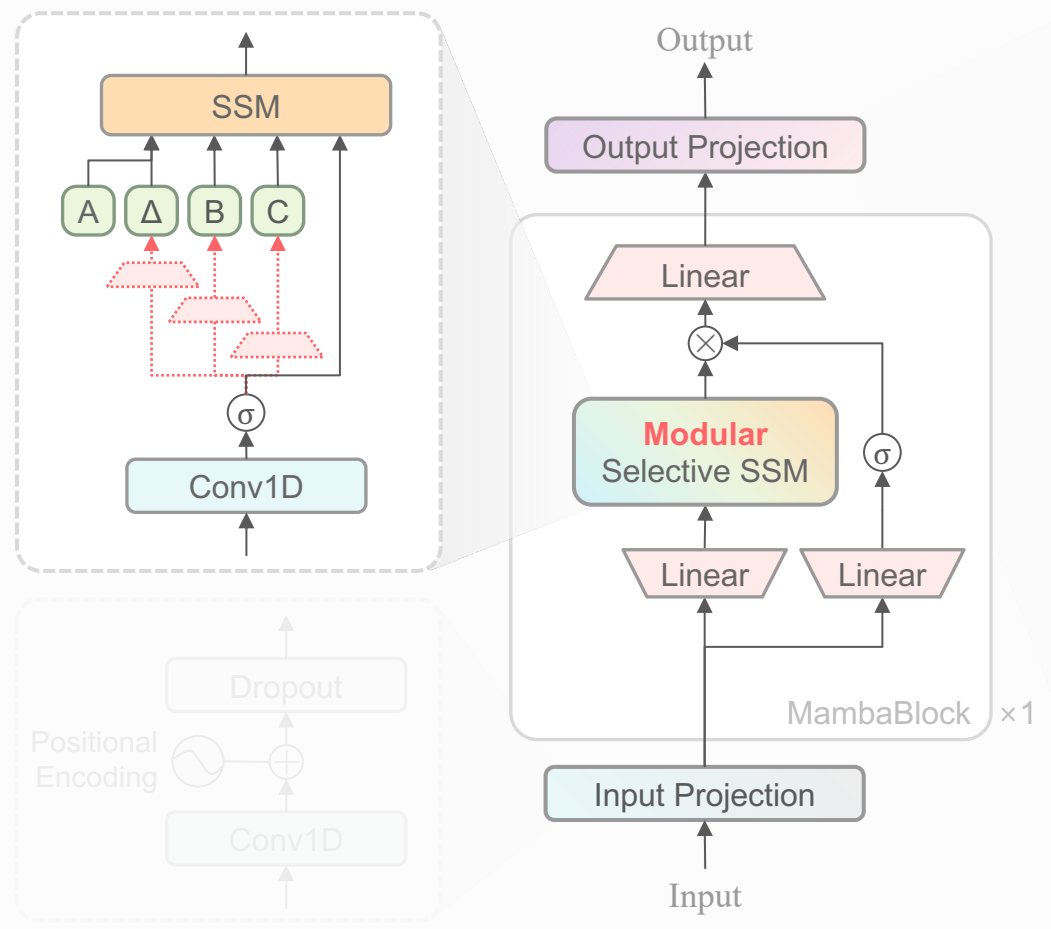


- decouple time variance of Δ , B , and C as a hyperparameter $\rightarrow \theta_{\Delta}, \theta_B, \theta_C \in \{0, 1\}$
- 8 combinations
 - if all $\theta = 0 \rightarrow$ LTI & channel-independent system
 - if all $\theta = 1 \rightarrow$ TV & channel-mixed system
- Ablation results showed that simpler, near time-invariant settings can be more effective.

σ : Learnable parameter
 (red dashed arrow) : Depends on time variability hyperparameters ($\theta_B, \theta_C, \theta_{\Delta}$)
 (green box) : Learnable parameter

Method

- (H3) Remove skip connection



- remove D from Mamba block

$$f_t^{(j)} = C_t s_t^{(j)} \quad (\text{no skip term}) \quad (12)$$

- Most time series classification models are in shallow layers (1 – 6 layers)

- InceptionTime (Chen et al., 2015) showed that the skip connection makes little difference across 85 UCR datasets..

σ : Learnable parameter
- - - : Depends on time variability
 : Learnable parameter

Method

- (H4) Aggregate via adaptive pooling

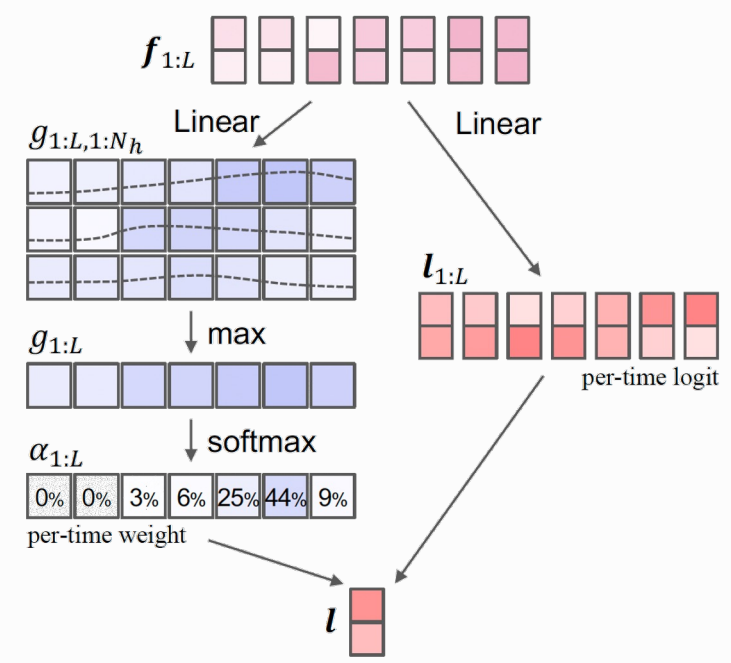
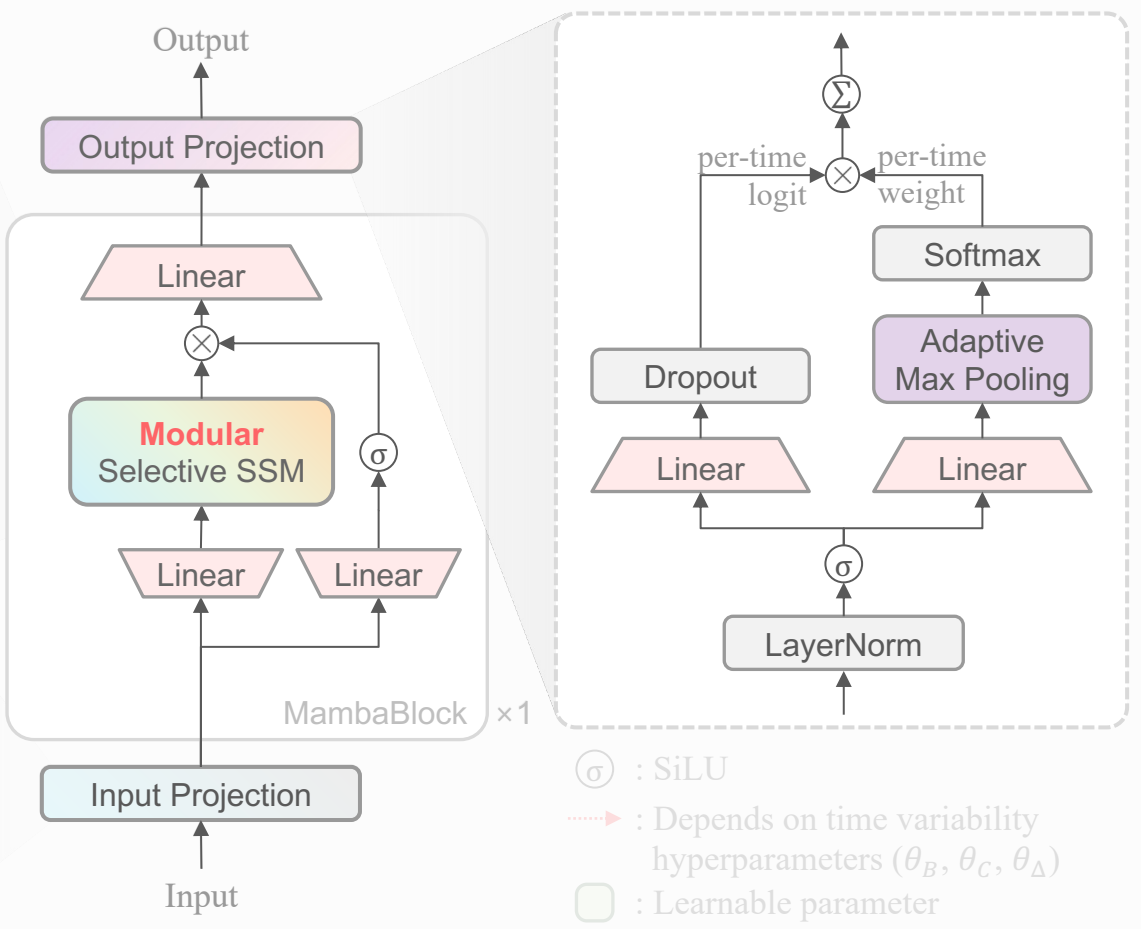
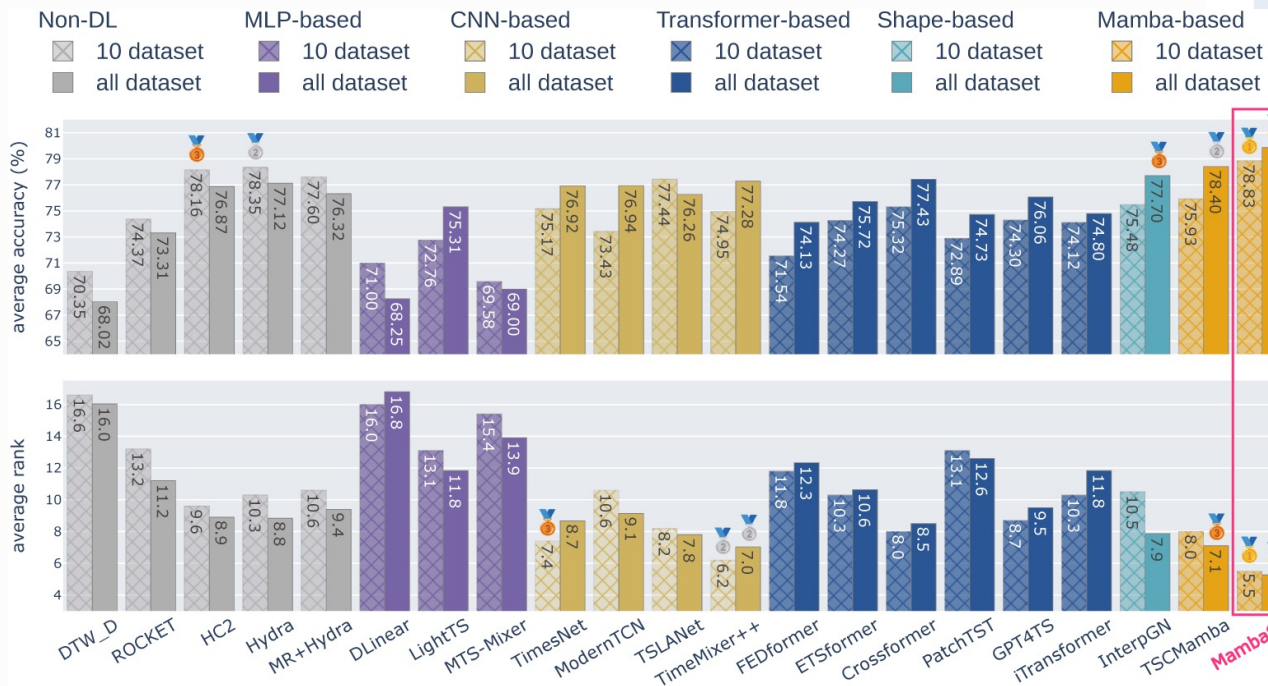
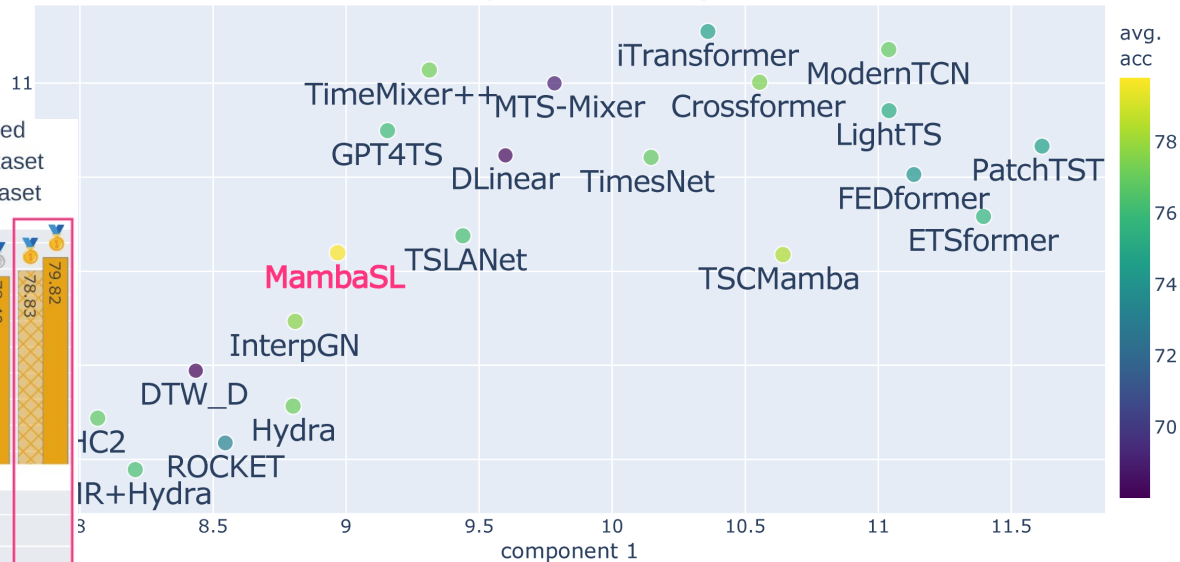


Figure 3: Illustration of proposed multi-head adaptive pooling with $(L, d_m, N_h, d_y) = (7, 2, 3, 2)$.

Results



Visualization of TSC models using UMAP (n_neighbors=12, metric=correlation)



- checkpoint_best
- LightTS.zip
 - DLinear.zip
 - ETSformer.zip
 - FEDformer.zip
 - GPT4TS.zip
 - InterpGN.zip
 - iTransformer.zip

- checkpoint_best
- MambaSL (ADFTD).zip
 - MambaSL (FLAAP).zip
 - MambaSL (inceptiontime-setting).zip
 - MambaSL (multilayer).zip
 - MambaSL (UCR, inceptiontime-setting).zip
 - MambaSL.zip

- checkpoint_best
- ModernTCN.zip
 - MTSMixer.zip
 - PatchTST.zip
 - TimeMixerPP.zip
 - TimesNet.zip
 - TSCMamba.zip
 - TSLANet.zip



ICLR



SNU IMLAB

Dept. of Industrial Engineering
Seoul National University

Thank you

Paper



Code



MambaSL : Exploring Single-Layer Mamba for Time Series Classification