

MAESTRO : Adaptive Sparse Attention and Robust Learning for Multimodal Dynamic Time Series

Payal Mohapatra Yueyuan Sui Akash Pandey Stephen Xia Qi Zhu
Northwestern University



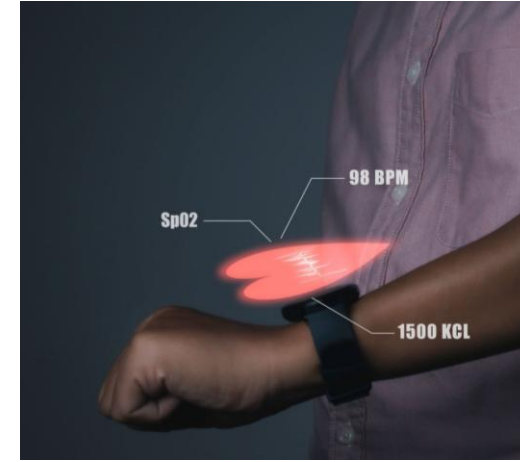
Time Series Data are Ubiquitous



Clinical Applications



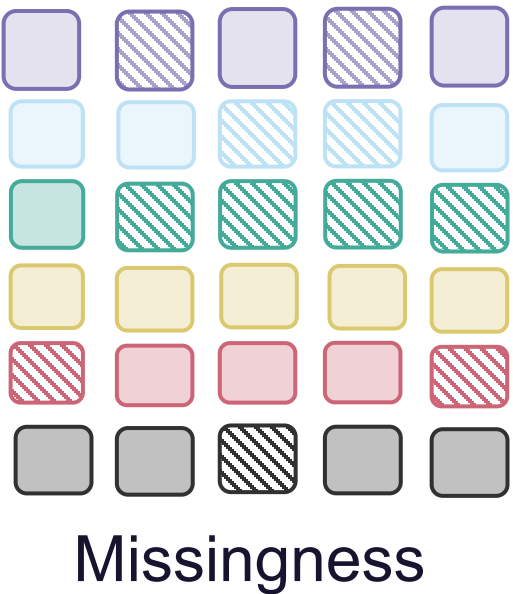
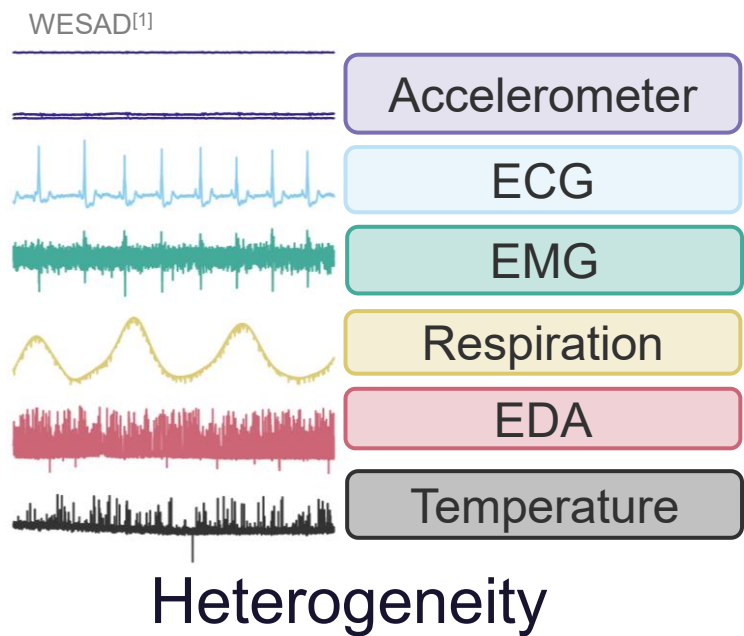
Environment Monitoring



Lifestyle Enhancement

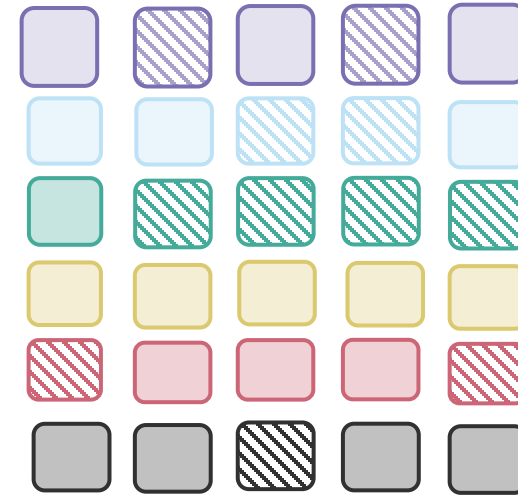
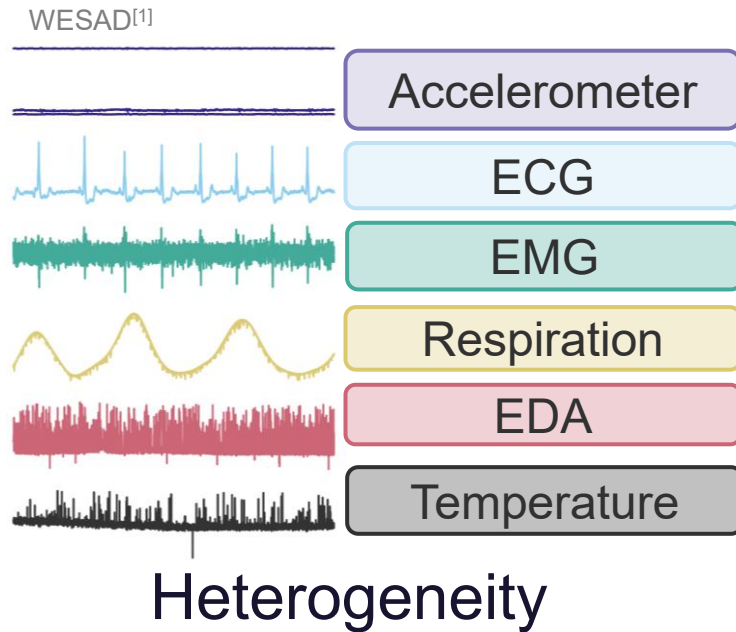


Challenges in Time Series Data Analyses



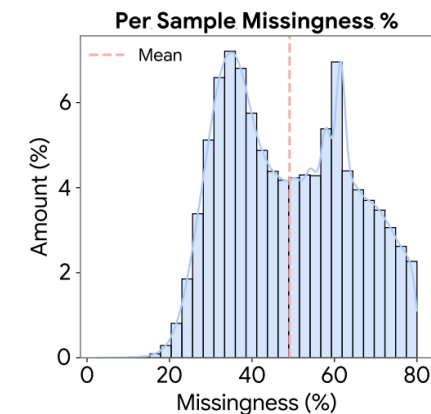
[1] Wearable Stress and Affect Detection, 2018

Challenges in Time Series Data Analyses



Missingness

40M hours of day-long multimodal sensor data from LSM-2^[2]



[1] Wearable Stress and Affect Detection, 2018

[2] LSM-2: Learning from Incomplete Wearable Sensor Data, 2025

Traditional Treatment of Time-series



Electrodermal Activity (EDA)



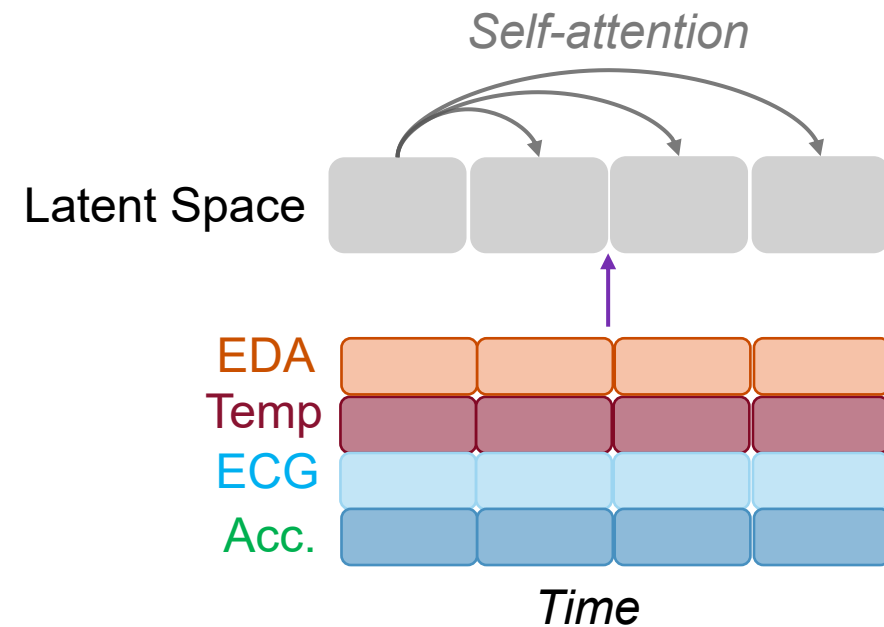
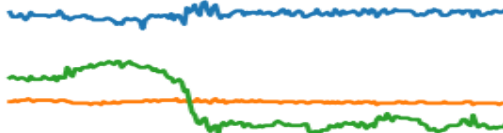
Skin Temperature (Temp)



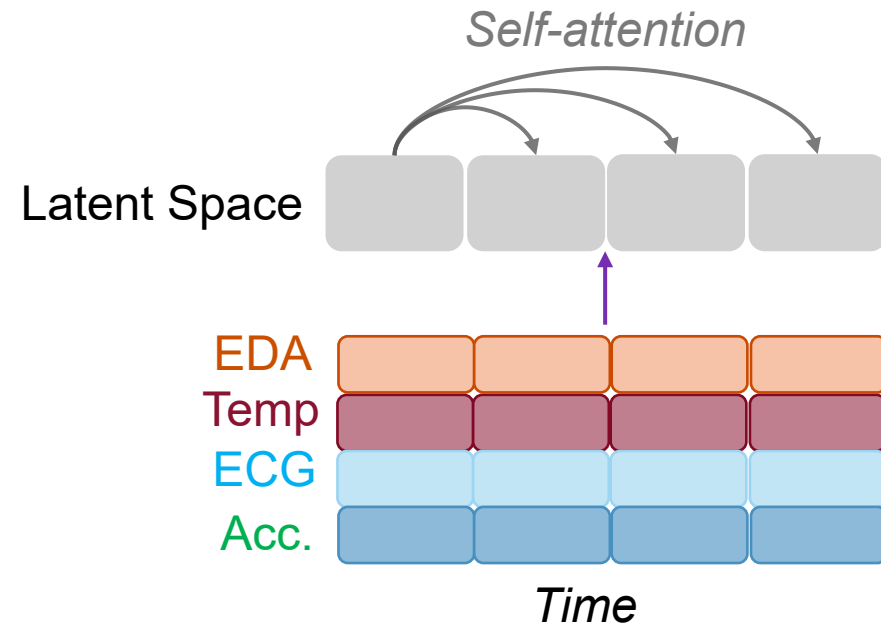
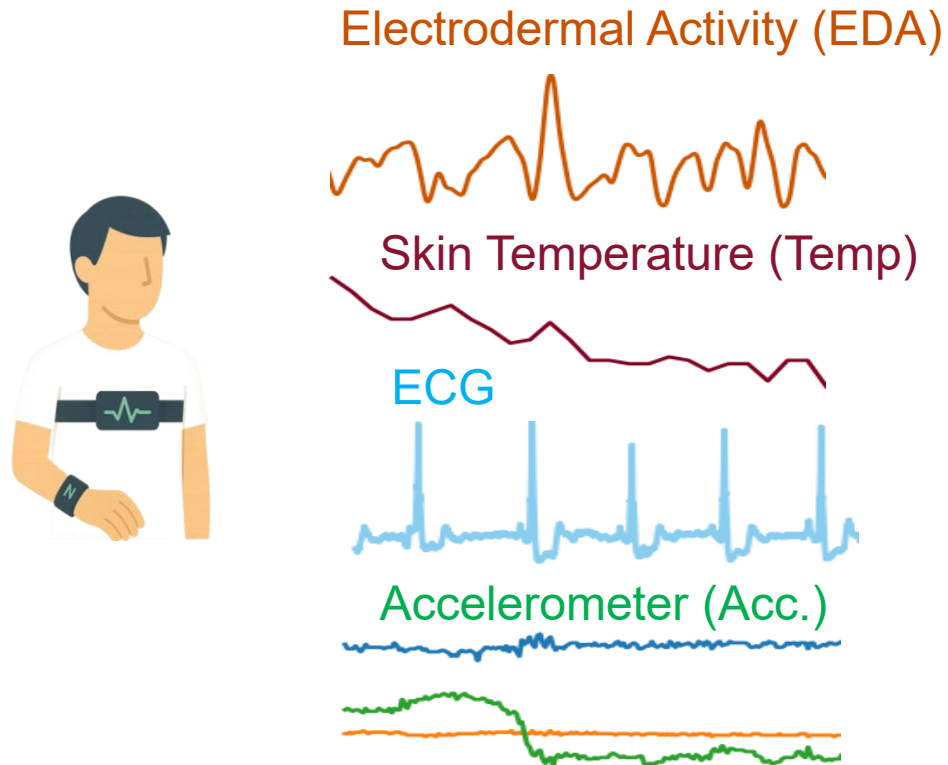
ECG



Accelerometer (Acc.)



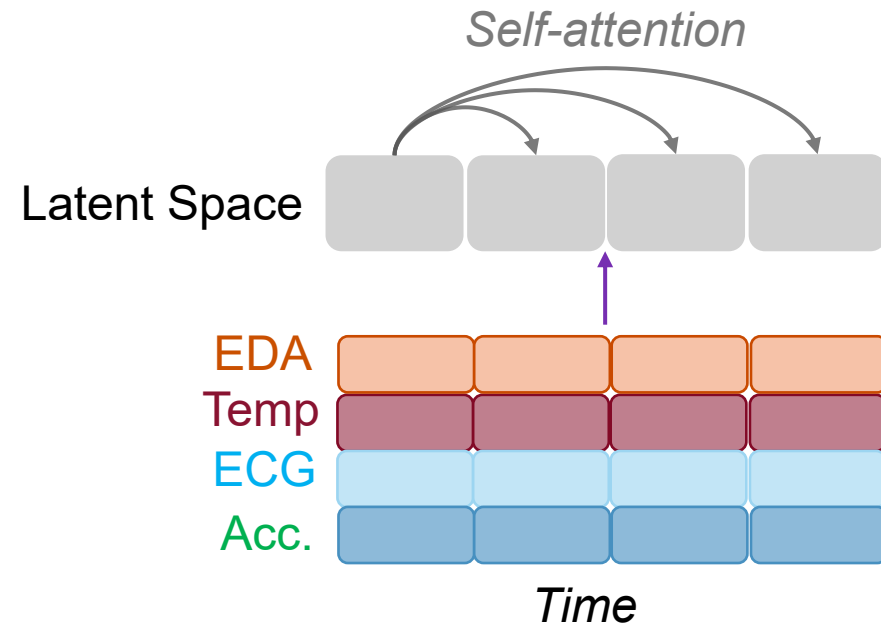
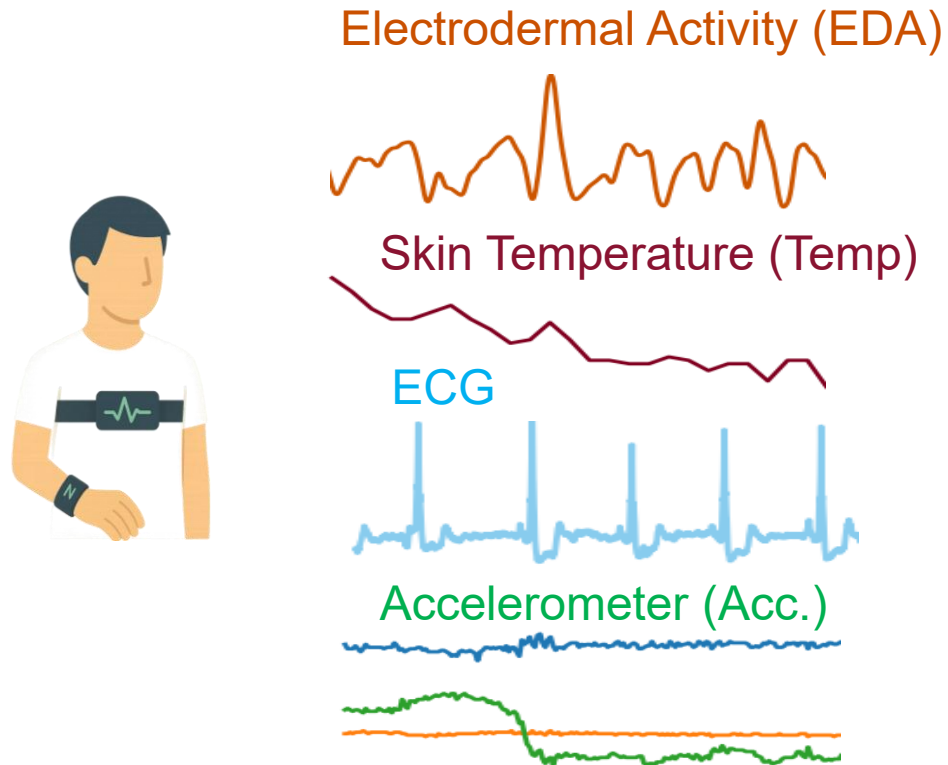
Traditional Treatment of Time-series



While effective in simple tasks :

- Does not model inherent heterogeneity
- Misses modeling inter-modal interactions effectively
- Cannot disentangle representation in case one modality is missing/corrupted.
- Cannot handle different sequence length across modalities.

Traditional Treatment of Time-series



Another approach is sensor fusion, but it is highly application-specific and often heuristic.

For example, in heart rate monitoring from wearables, several proposals have been made in the literature, such as:

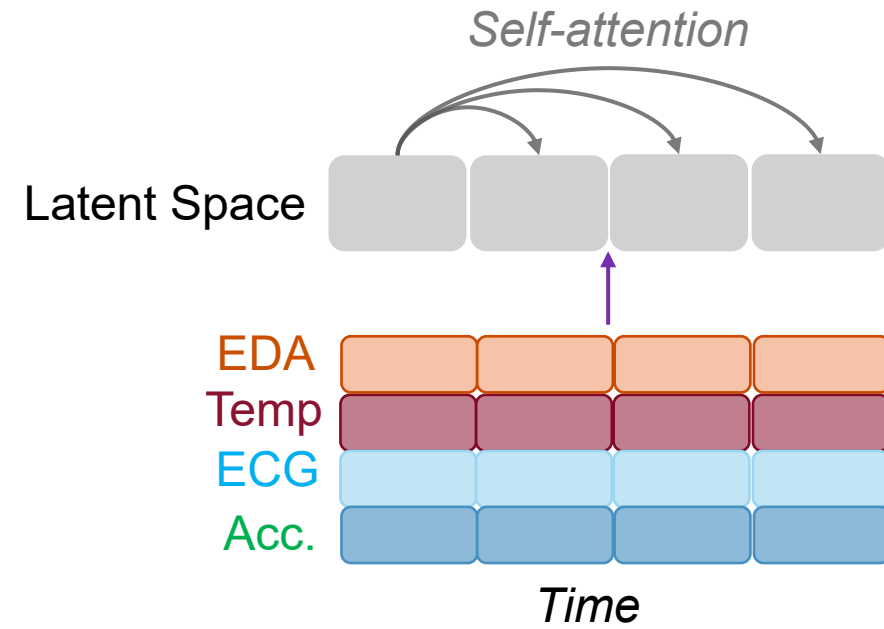
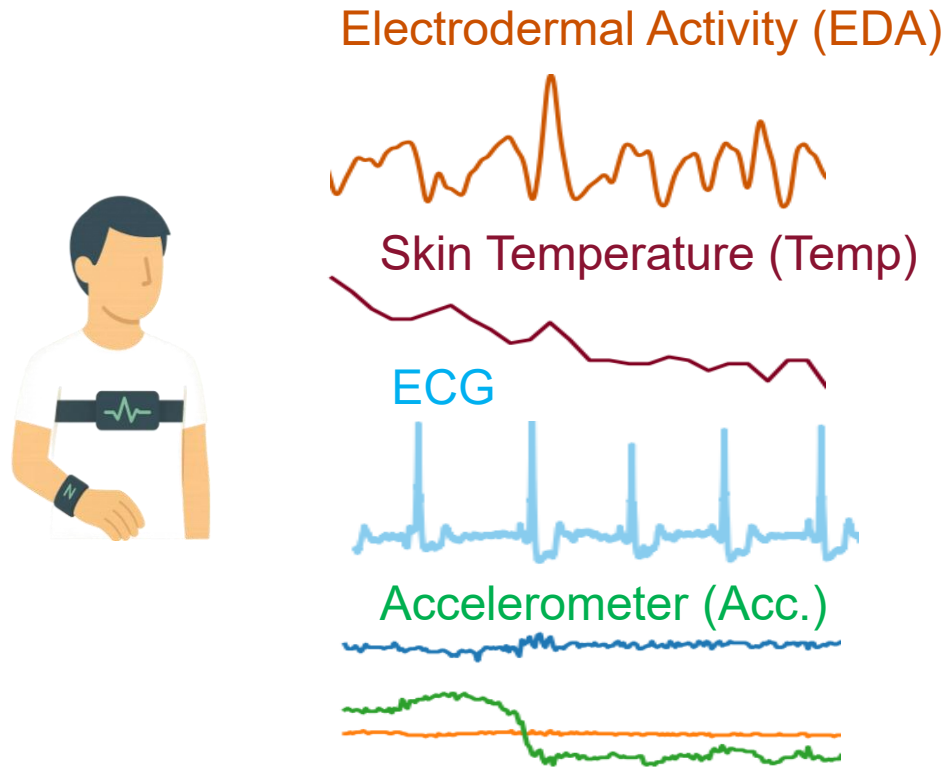
- *early fusion using multi-wavelength photoplethysmography (PPG)¹,*
- *multi-site PPG², and*
- *late fusion with temperature³.*

[1] Meier & Holz (2024). *Effect of wavelength on PPG reliability outdoors* — CISS.

[2] Meier & Holz (2024). *PPG accuracy across body locations and motion* — CISS.

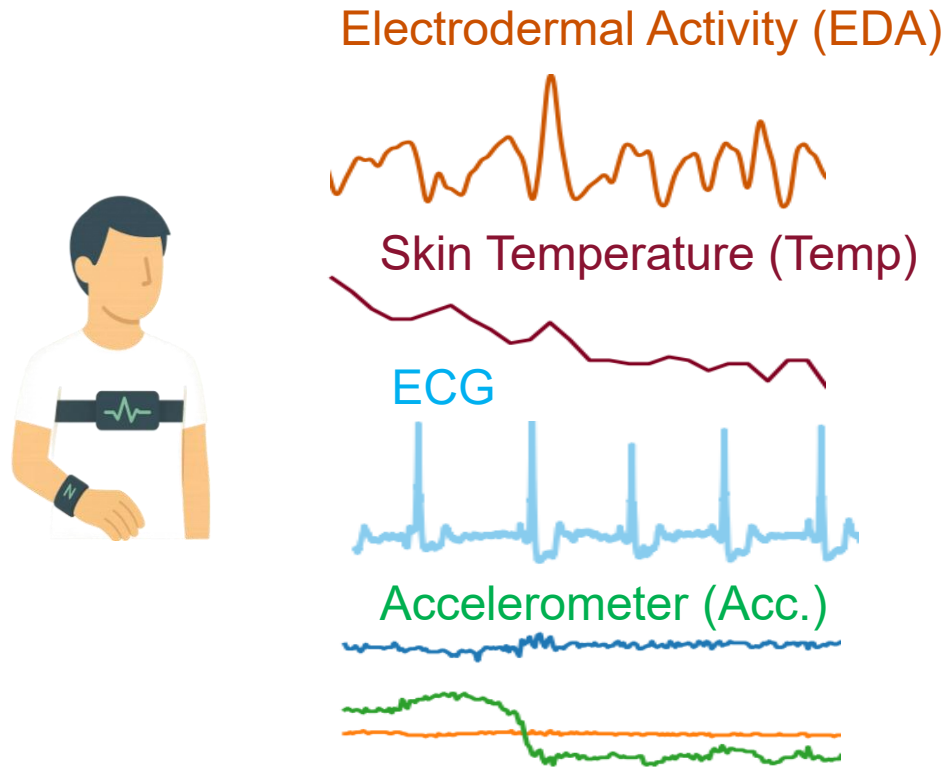
[3] Meier, Demirel & Holz (2024). *WildPPG: Real-world long-duration PPG dataset* — NeurIPS.

Traditional Treatment of Time-series



Need to view as multimodal data
instead of multivariate-time-series!

Traditional Treatment of Time-series



Generally, video, audio, and text representations are considered multimodal. However, even though time-series data are represented in the same numerical form, they can be highly heterogeneous.

Need to view as multimodal data instead of multivariate-time-series!

Considerations for Multimodal Real-world Time-series

➔ Need to identify primary modality^[1].

➔ Assumption of high mutual information among modalities^[2, 3, 4].

➔ Pairwise interaction modeling^[5, 6].

🔒 *Apriori* of primary modality is not always guaranteed.

🔒 Heterogenous modalities.

🔒 Number of modalities can be greater than 10. Combinatorially expansive!

🔒 Random missingness due to sensor malfunction.

[1] IMAGEBIND, CVPR 2023 (highlight paper)

[2] VATT, Neurips 2021

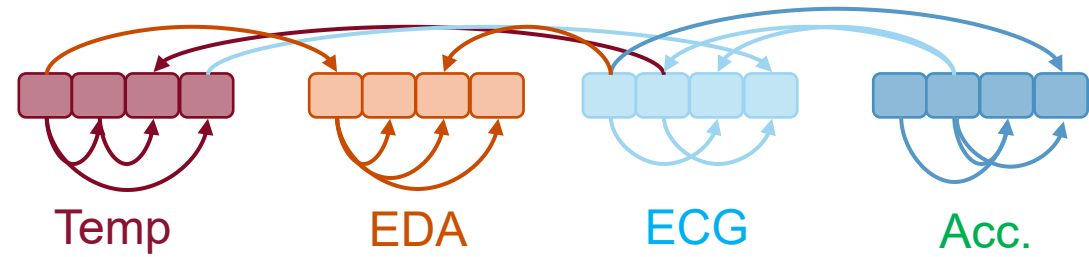
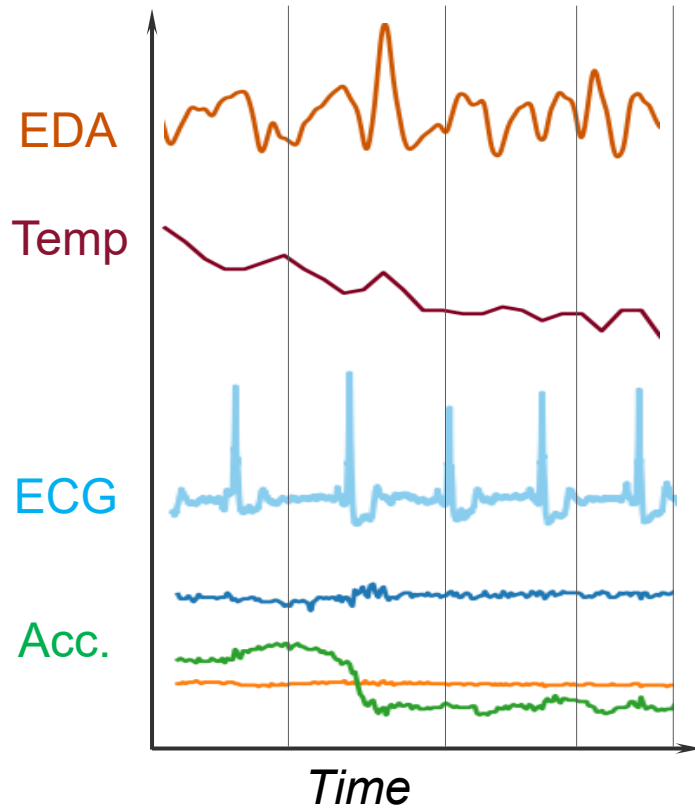
[3] Factorized Contrastive Learning, Neurips 2023

[4] Multimodal Fusion Interactions, ICMI 2023

[5] MULT, ACL 2020

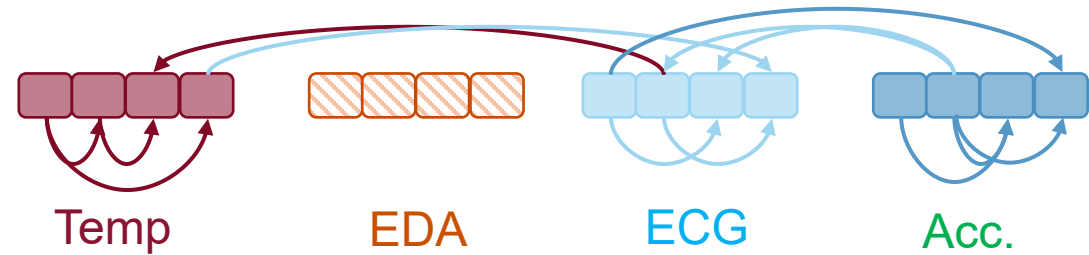
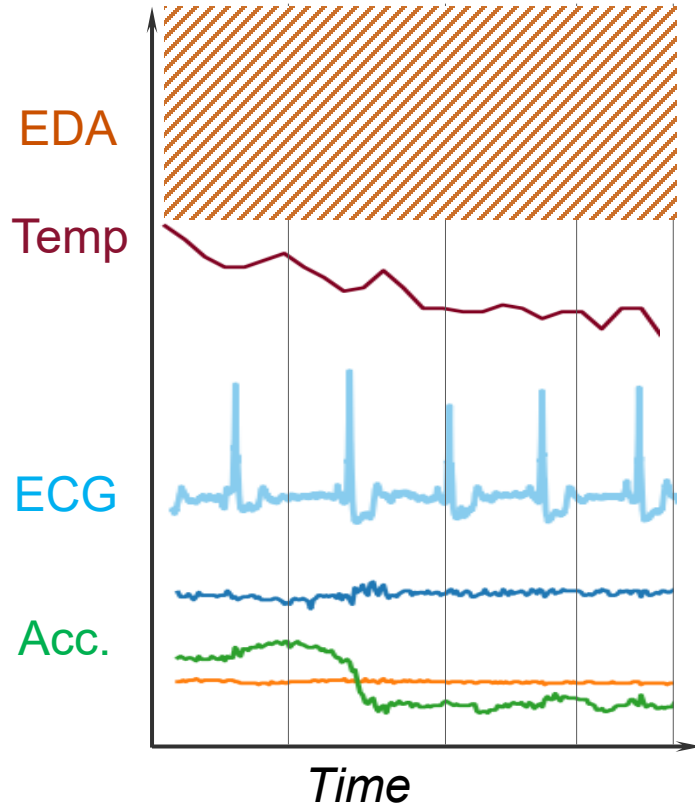
[6] MMOE, EMNLP 2024

Cross-modal-attention for Multimodal Time-series



Cross-attention can allow learning task-relevant modality interaction.

Cross-modal-attention for Multimodal Time-series



Cross-attention can allow learning from arbitrary modality combinations.

But applying Cross-modal-attention through Long Multimodal Time-series increases the computational complexity!

Canonical Self-Attention^[1]

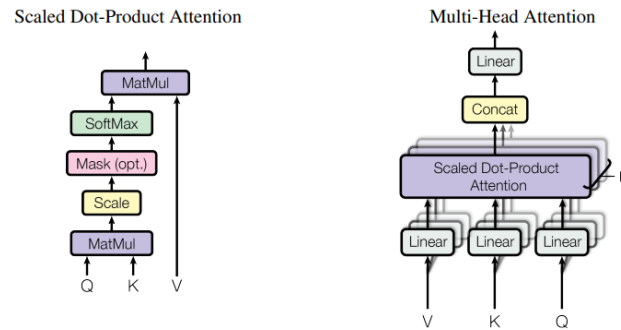


Figure 2: (left) Scaled Dot-Product Attention. (right) Multi-Head Attention consists of several attention layers running in parallel.

$$\mathcal{A}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}} \right) \mathbf{V}$$

Point-wise self-attention for a sequence length of L , has a quadratic computational complexity $\rightarrow \mathcal{O}(L^2)$

Consider M modalities each of sequence length L , then the computational complexity increases,

$$\mathcal{O}(M^2 L^2)$$

But applying Cross-modal-attention through Long Multimodal Time-series increases the computational complexity!

Canonical Self-Attention^[1]

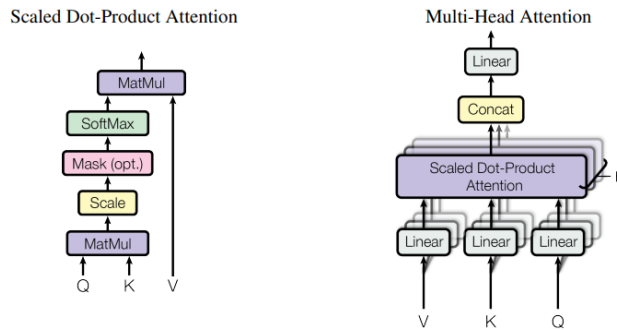


Figure 2: (left) Scaled Dot-Product Attention. (right) Multi-Head Attention consists of several attention layers running in parallel.

$$\mathcal{A}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}} \right) \mathbf{V}$$

Point-wise self-attention for a sequence length of L , has a quadratic computational complexity $\rightarrow \mathcal{O}(L^2)$

Consider M modalities each of sequence length L , then the computational complexity increases,

$$\mathcal{O}(M^2 L^2)$$

Sparse
Attention

$$\mathcal{O}(ML \log(ML))$$

Handling long time-series through sparse attention (Overview)

- Point-wise self-attention for a sequence length of L , has a quadratic computational complexity.

$$\mathcal{A}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}} \right) \mathbf{V}$$

ProbSparse Attention - $\mathcal{A}_s(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax} \left(\frac{\bar{\mathbf{Q}}\mathbf{K}^\top}{\sqrt{d}} \right) \mathbf{V}$

- Stacking N encoder layers further increases the memory consumption and computational complexity.

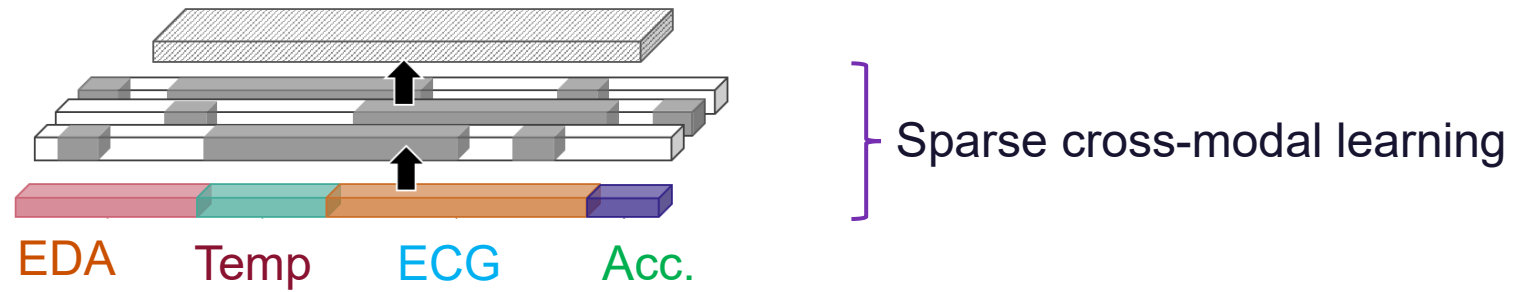
$$\hat{s} = s + \text{PE}_{\text{sin}}(s)$$

$$\bar{s} = \mathcal{A}_s(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax} \left(\frac{\bar{\mathbf{Q}}\mathbf{K}^\top}{\sqrt{d}} \right) \mathbf{V}$$

$$\dot{s} = \bar{s} + \hat{s}$$

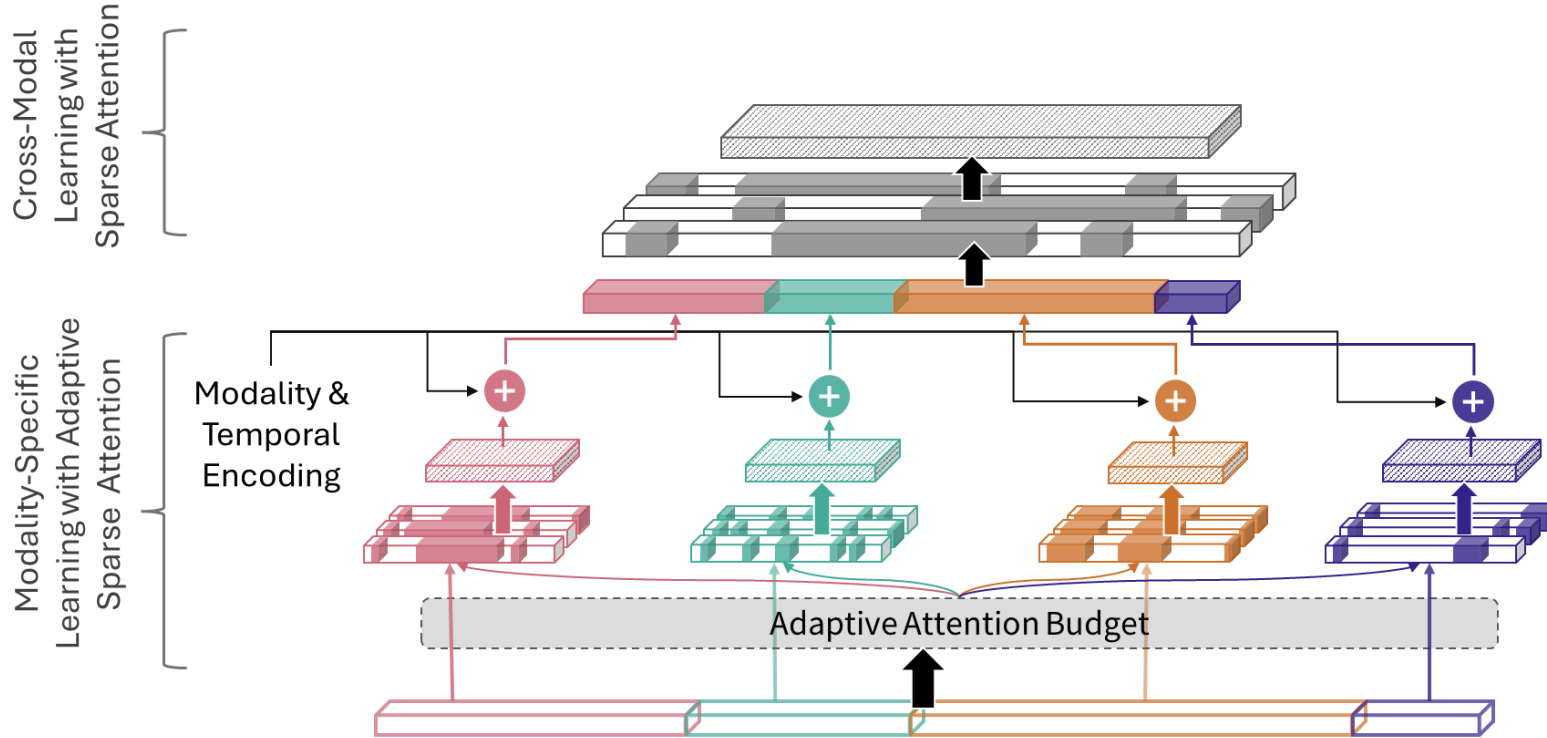
$$z = \text{distil}(\dot{s}) + \text{maxpool}(\hat{s})$$

Handling long multimodal time-series through sparse attention



1. Adaptive Attention Budget per modality

Top- \mathbf{v} queries in ProbSparse Computation : \mathbf{v} is modulated by modality's relevance and availability



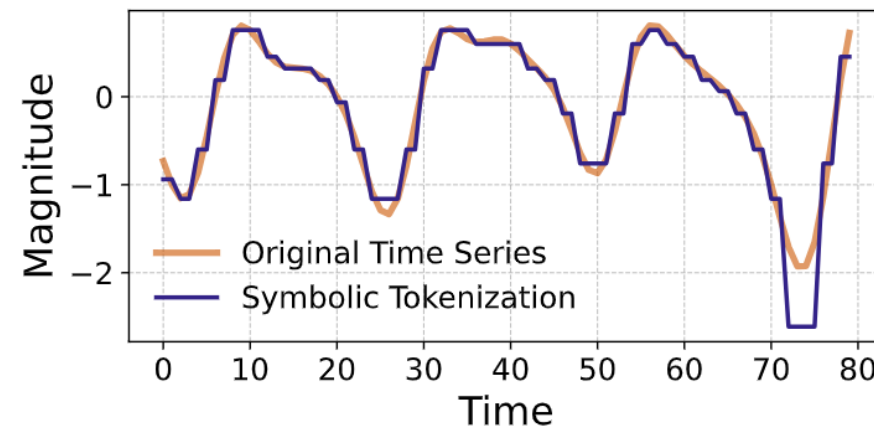
2. Symbolic Tokenization

1. Converts time-series to discrete *symbols*^[1].

- Has some nice properties – guarantees a lower bound Euclidean distance between the symbolic time-series and the original time-series.
- We extend it under some assumptions that this tokenization preserves multimodal relational structure.

2. We can reserve a symbol for *missing* data naturally.

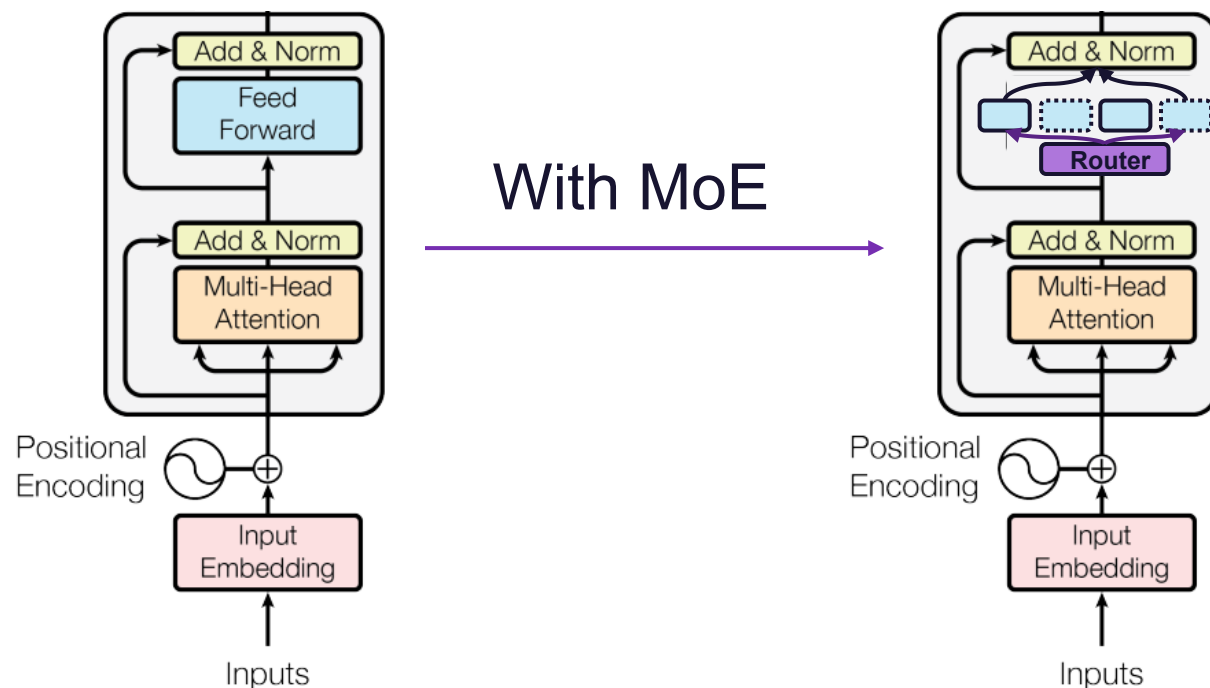
3. We can compress the signal further reducing the sequence length.



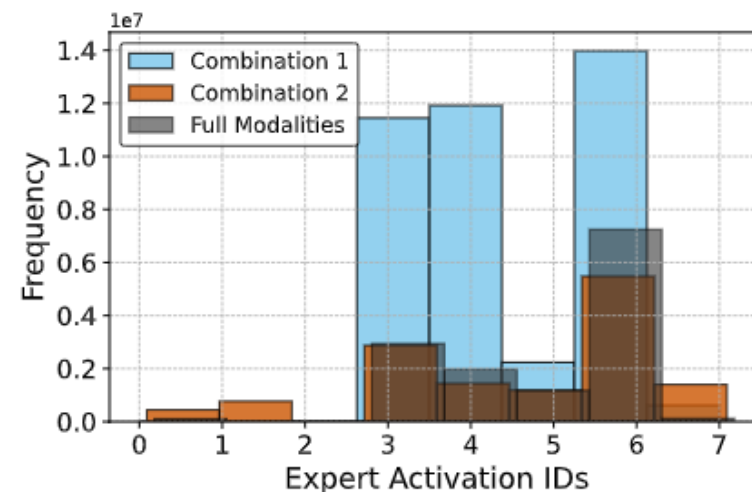
[1] Experiencing SAX: a Novel Symbolic Representation of Time Series, Data Mining and Knowledge Discovery, 2007

3. Handling Missingness through Mixture-of-Experts

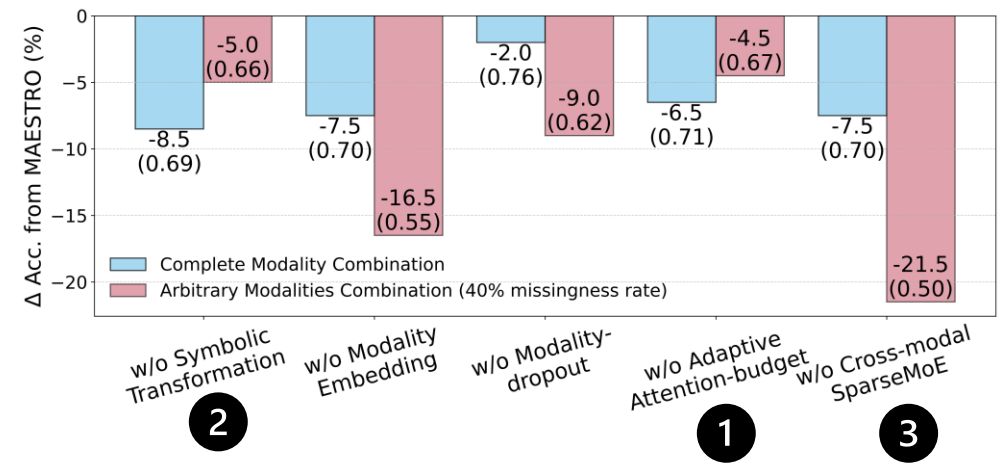
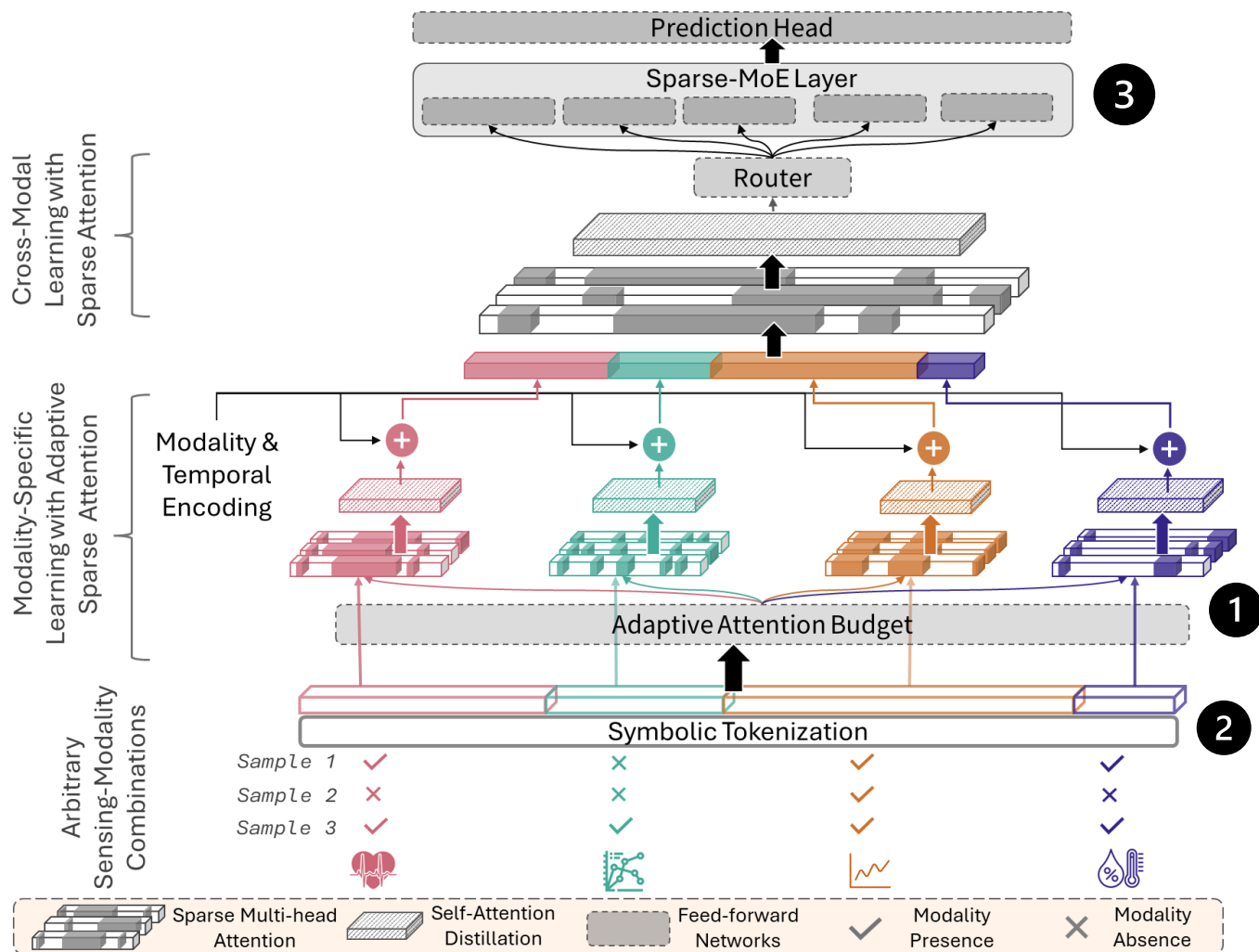
Instead of the fixed Feed-forward layer of the transformer, we can use a mixture of experts(MoE) for implicit modality specialization.



Vanilla Transformer^[1]



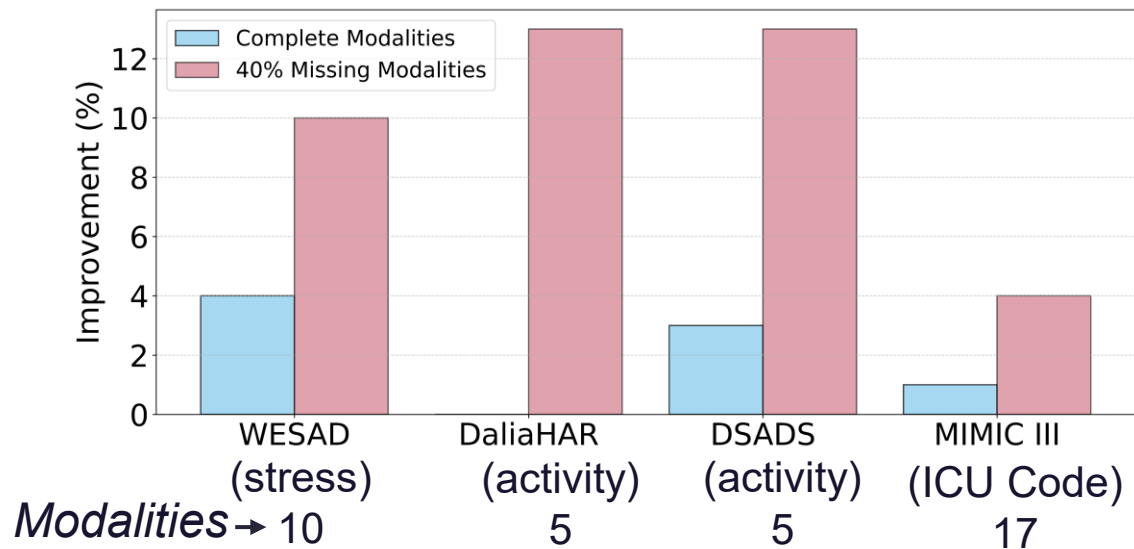
Overall MAESTRO Framework



Ablation of Key Components

Key Results : MAESTRO

Average of 4% performance improvement with complete and 8% with arbitrary set of modalities.



Computational Efficiency

Model	Acc. \uparrow	MMAC \downarrow	GFLOPs \downarrow	Params (M)
<i>Multivariate Models</i>				
iTransformer	0.67 ± 0.05	2833	5.73	12.82
Transformer	0.63 ± 0.02	4331	8.66	1.68
<i>Multimodal Models</i>				
FuseMoE	0.47 ± 0.41	6524	13.05	0.67
MULT	0.60 ± 0.42	13324	26.65	3.71
ShaSpec	0.62 ± 0.51	4556	9.11	216
MAESTRO	0.77 ± 0.04	3066	6.13	1.39
– Full-Attn (Per-Modal)	0.80 ± 0.03	3769	7.54	1.40
– Full-Attn (Cross-Modal)	0.77 ± 0.07	3496	6.99	1.39
– All Full-Attention	0.75 ± 0.05	4205	8.42	1.39
– All Full-Attention (no MoE)	0.78 ± 0.04	4392	8.78	1.39

Thank you.