# A   More Details on Time Series Line Graph Image Creation

**Implementation.**  The time-series-to-image transformation can be implemented using the Matplotlib package[1] with the following few lines of code.

```python
def TS2Image(t, v, D, colors, image_height, image_width, grid_height, grid_width):
    import matplotlib.pyplot as plt
    plt.figure(figsize=(image_height/100, image_width/100), dpi=100)
    for d in range(D): # enumerate the multiple variables
        plt.subplot(grid_height, grid_width, d+1) # position in the grid
        # plot line graph of variable d
        plt.plot(t[d], v[d], color=colors[d], linestyle="-", marker="*")
```

In addition to the designs mentioned in the main paper for plotting line graph images, we also explore the following aspects.

**Axis Limits of Line Graphs.**  The axis limits determine the plot area of the line graphs and the range of displayed timestamps and values. By default, we set the limits of the x-axis and y-axis as the ranges of all the observed timestamps and values across the dataset. However, we found that some extreme observed values for some variables can make the range of the y-axis very large, causing most plotted points of observations to cluster in a small area and resulting in flat line graphs. Common normalization and standardization methods will not solve this issue, as the relative magnitudes remain unchanged in the created images. We thus tried the following strategies to remove extreme values and narrow the range of the y-axis:

- Interquartile Range (IQR): IQR is one of the most extensively used methods for outlier detection and removal. The interquartile range is calculated based on the first and third quartiles of all the observed values of each variable in the dataset and then used to calculate the upper and lower limits.

- Standard Deviation (SD): The upper and lower boundaries are calculated by taking 3 standard deviations from the mean of observed values for each variable across the dataset. This method usually assumes the data is normally distributed.

- Modified Z-score (MZ): A z-score measures how many standard deviations away a value is from the mean and is similar to the standard deviation method to detect outliers. However, z-scores can be influenced by extreme values, which modified z-scores can better handle. We set the upper and lower limits as the values whose modified z-scores are 3.5 and -3.5.

We show examples of the created images with these strategies in Figure 1.

Table 1: Ablation study on different strategies to decide the line graph limit. The default strategy is to directly set the axis limit as the range of all observed values on the dataset. "IQR", "SD", and "MZS' denote three strategies to remove extreme value, *i.e.,* Interqurtile Range, Standard Deviation, and Modified Z-score. The reported numbers are averaged on 5 data splits.

| Strategies | P19 | | P12 | | PAM | | | |
|---|---|---|---|---|---|---|---|---|
| | AUROC | AUPRC | AUROC | AUPRC | Accuracy | Precision | Recall | F1 score |
| Default | $\mathbf{89.4} \pm 1.9$ | $\mathbf{52.8} \pm 3.8$ | $\mathbf{85.6} \pm 1.1$ | $\mathbf{49.8} \pm 2.5$ | $96.1 \pm 0.7$ | $96.8 \pm 1.1$ | $96.5 \pm 0.7$ | $96.6 \pm 0.9$ |
| IQR | $88.2 \pm 0.8$ | $49.6 \pm 1.7$ | $84.5 \pm 1.1$ | $48.9 \pm 2.6$ | $95.9 \pm 0.7$ | $96.8 \pm 0.7$ | $96.1 \pm 0.7$ | $96.4 \pm 0.7$ |
| SD | $87.4 \pm 1.6$ | $51.2 \pm 3.6$ | $84.6 \pm 1.7$ | $47.1 \pm 2.9$ | $\mathbf{96.6} \pm 0.9$ | $\mathbf{97.1} \pm 0.8$ | $\mathbf{97.0} \pm 0.6$ | $\mathbf{97.0} \pm 0.7$ |
| MZS | $87.3 \pm 1.0$ | $50.8 \pm 3.7$ | $84.3 \pm 1.4$ | $47.1 \pm 2.1$ | $96.0 \pm 1.1$ | $96.8 \pm 0.99$ | $96.4 \pm 0.9$ | $96.6 \pm 0.9$ |

The performance comparison of models trained on images created with different strategies is shown in Table 1. We observe that the methods that remove extreme values hurt the performance, except for SD on the PAM dataset. Although these methods narrow the value range and highlight the dynamic patterns of line graphs, they discard the extreme values which might be informative themselves. This observation suggests that our approach may not require additional data preprocessing on the time series, further demonstrating its advantage in simplicity.
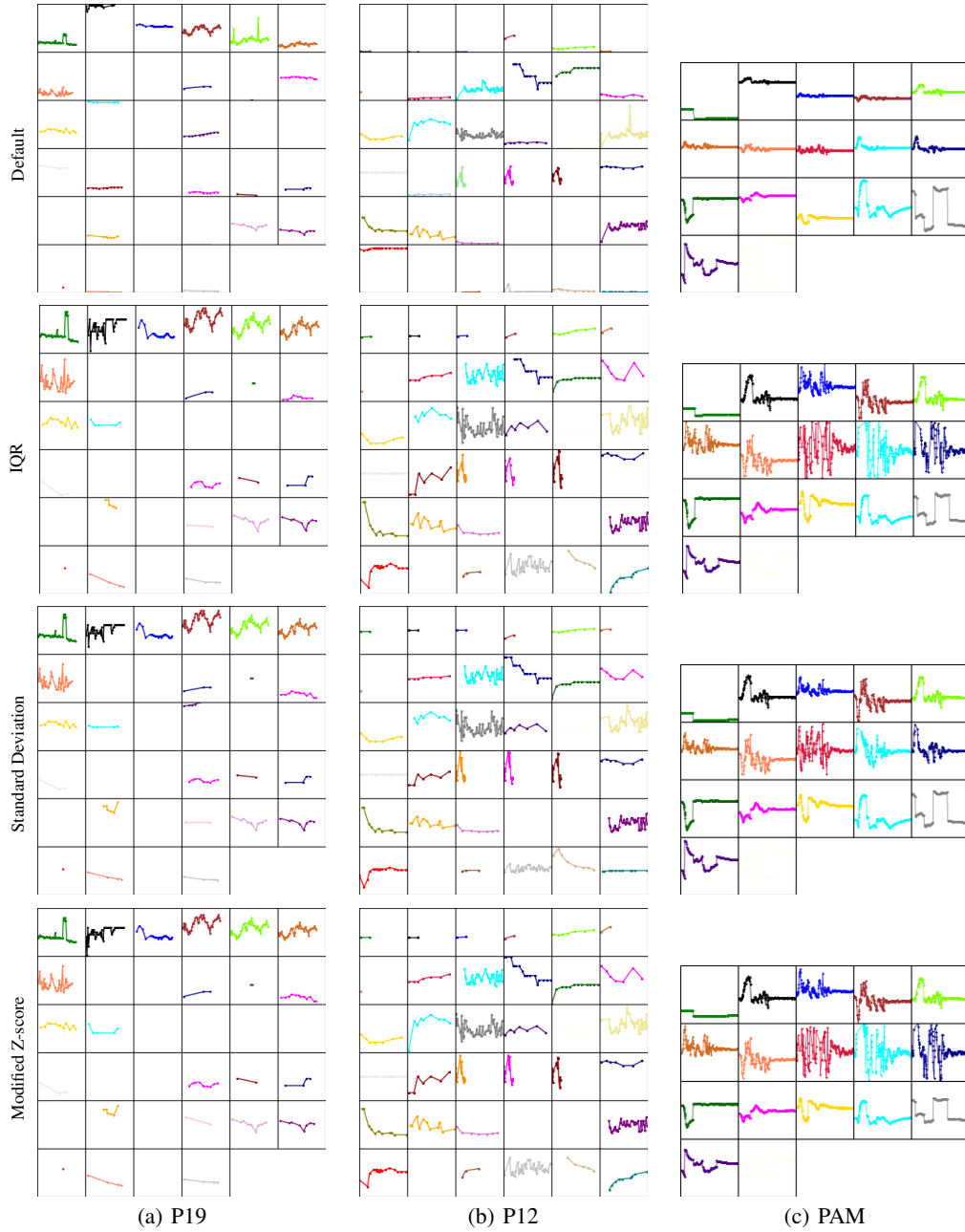
---
[1]https://matplotlib.org/

1

Figure 1: The images created with different strategies for three samples from P19, P12, and PAM dataset, respectively (sample "p000019" for P19, "132548" for P12, and "0" for PAM).

Table 2: Ablation study on grid layouts and image sizes on P19.

| Grid Layout | Image Size | AUROC | AUPRC |
|---|---|---|---|
| $4 \times 9$ | $256 \times 576$ | $87.4_{\pm 1.9}$ | $48.1_{\pm 4.5}$ |
| $5 \times 7$ | $320 \times 448$ | $87.9_{\pm 1.9}$ | $49.6_{\pm 2.7}$ |
| $6 \times 6$ | $384 \times 384$ | $\mathbf{89.4}_{\pm 1.9}$ | $\mathbf{52.8}_{\pm 3.8}$ |
| $6 \times 6$ | $224 \times 224$ | $88.7_{\pm 1.4}$ | $52.3_{\pm 0.6}$ |

Table 3: Ablation study on grid layouts and image sizes on P12.

| Grid Layout | Image Size | AUROC | AUPRC |
|---|---|---|---|
| $4 \times 9$ | $256 \times 576$ | $84.0_{\pm 1.4}$ | $47.9_{\pm 2.6}$ |
| $5 \times 8$ | $320 \times 512$ | $84.1_{\pm 1.6}$ | $47.2_{\pm 2.3}$ |
| $6 \times 6$ | $384 \times 384$ | $\mathbf{85.6}_{\pm 1.1}$ | $\mathbf{49.8}_{\pm 2.5}$ |
| $6 \times 6$ | $224 \times 224$ | $85.2_{\pm 2.1}$ | $48.8_{\pm 3.7}$ |

Table 4: Ablation study on grid layouts and image sizes on PAM.

| Grid Layout | Image Size | Accuracy | Precision | Recall | F1 score |
|---|---|---|---|---|---|
| $2 \times 9$ | $128 \times 576$ | $95.9_{\pm 1.4}$ | $96.5_{\pm 1.0}$ | $\mathbf{95.9}_{\pm 1.2}$ | $96.0_{\pm 0.5}$ |
| $3 \times 6$ | $192 \times 384$ | $\mathbf{96.1}_{\pm 0.8}$ | $96.7_{\pm 0.5}$ | $95.9_{\pm 0.9}$ | $96.2_{\pm 0.7}$ |
| $4 \times 5$ | $256 \times 320$ | $\mathbf{96.1}_{\pm 0.7}$ | $\mathbf{96.8}_{\pm 1.1}$ | $96.5_{\pm 0.7}$ | $\mathbf{96.6}_{\pm 0.9}$ |
| $4 \times 5$ | $224 \times 224$ | $95.9_{\pm 0.6}$ | $96.7_{\pm 0.8}$ | $\mathbf{95.9}_{\pm 0.6}$ | $96.3_{\pm 0.7}$ |

**Grid Layout and Image Size.** We conducted experiments to study the impact of grid layouts and image sizes on the performance of our approach. For a fair comparison of different grid layouts, we fixed the size of each grid cell as $64 \times 64$ and altered the grid layouts. The results on the P19, P12, and PAM datasets are listed in Table 2, Table 3, and Table 4, respectively. We observed that the square grid layouts consistently produced good results on all three datasets. We conjecture that this is because the square layout ensures that the distance between any two line graphs is shortest. We also tested the performance with the standard image size of $224 \times 224$ and found that the differences were marginal, indicating the robustness of our approach to various image sizes.

# B   More Experimental Details

## B.1   Datasets

We used the datasets processed by [7], whose details are given below.

**P19: PhysioNet Sepsis Early Prediction Challenge 2019.** [2] The P19 dataset [5] consists of clinical data for 38,803 patients, and aims to predict whether sepsis will occur within the next 6 hours. The dataset includes 34 irregularly sampled sensors with 8 vital signs and 26 laboratory values for each patient, as well as 6 demographic features. To process the static features, we use templates outlined in Table 5, and utilize a pre-trained Roberta-base model to extract textual features. These textual features are then combined with visual features obtained from the vision transformer to perform binary classification. The dataset is highly imbalanced with only 4% of samples being positive, and has a missing ratio of 94.9%.

**P12: PhysioNet Mortality Prediction Challenge 2012.** [3] P12 dataset [2] includes clinical data from 11,988 ICU patients, with 36 irregularly sampled sensor observations and 6 static demographic features provided for each patient. The goal is to predict patient mortality, which is a binary classification task. The dataset is highly imbalanced, with around 86% of samples being negative. The missing ratio of the dataset is 88.4%.

**PAM: PAMAP2 Physical Activity Monitoring.** [4] The PAM dataset originally contains data of 18 physical activities with 9 subjects wearing 3 inertial measurement units. However, to make it suitable for irregular time series classification, [7] excluded the ninth subject due to its short length of sensor readouts, and 10 out of the 18 activities that had less than 500 samples were also excluded. As a result, the task is an 8-way classification with 5,333 samples, each with 600 continuous observations. To simulate the irregular time series setting, 60% of the observations are randomly removed. No static features are provided, and the 8 categories are approximately balanced. The missing ratio is 60.0%.

---

[2]https://physionet.org/content/challenge-2019/1.0.0/
[3]https://physionet.org/content/challenge-2012/1.0.0/
[4]https://archive.ics.uci.edu/ml/datasets/pamap2+physical+activity+monitoring

Table 5: Templates for transforming static features to natural language sentences.

| Dataset | Static features | Template | Example |
|---------|-----------------|----------|---------|
| P19 | *Age*, *Gender*, *Unit1 (medical ICU)*, *Unit2 (surgery ICU)*, *HospAdmTime*; *ICULOS (ICU length-of-stay)* | A patient is {*Age*} years old, {*Gender*}, went to {*Unit1&Unit2*} {*HospAdmTime*} hours after hospital admit, had stayed there for {*ICULOS*} hours. | A patient is 65 years old, female, went to the medical ICU 10 hours after hospital admit, had stayed there for 20 hours. |
| P12 | *RecordID*, *Age*, *Gender*, *Height* (cm), *ICUType*, *Weight* (kg) | A patient is {*Age*} years old, {*Gender*}, {*Height*} cm, {*Weight*} kg, stayed in {*ICUType*}. | A patient is 48 years old, male, 171 cm, 78 kg, stayed in surgical ICU. |

Table 6: Ablation studies on different methods to encode static features.

| Methods | P19 | | P12 | |
|---------|-----|-----|-----|-----|
| | AUROC | AUPRC | AUROC | AUPRC |
| Raindrop | $87.0 \pm 2.3$ | $51.8 \pm 5.5$ | $82.8 \pm 1.7$ | $44.0 \pm 3.0$ |
| Swin | $89.4 \pm 1.8$ | $50.2 \pm 3.0$ | $84.3 \pm 0.6$ | $49.3 \pm 3.7$ |
| Swin-MLP | $88.6 \pm 1.3$ | $51.4 \pm 3.7$ | $84.6 \pm 0.9$ | $48.7 \pm 3.2$ |
| Swin-Roberta | $\mathbf{89.4} \pm 1.9$ | $\mathbf{52.8} \pm 3.8$ | $85.6 \pm 1.1$ | $\mathbf{49.8} \pm 2.5$ |

## B.2 Experiments on Static Features

Time series data is often associated with information from other modalities, such as the textual clinical notes in electronic health records (EHRs) in the healthcare domain. Our approach is naturally suitable for incorporating such information since we convert time series data to images, and thus various vision-language and multi-modal techniques can be utilized to incorporate the visual (time series) information and information from other modalities. For example, the CLIP [4] learns a shared hidden feature space where the paired image and text stay close. Under our framework, such a shared space can also be learned for the paired visual time series images and textual clinical notes, which is our future direction. It also paves the way for the application of multi-modal models such as GPT-4 [3] to handle the visualized time series data and the clinical notes simultaneously. In our current experiments, we used a text encoder, Roberta-base, to encode textual demographic information in the P19 and P12 datasets. We also experimented with normalizing the original categorical features and encoding them using an MLP as in previous work, and compare with the strong baseline, Raindrop. The results are shown in Table 6. We observe that even without using static features, our method has already outperformed Raindrop. In addition, utilizing Roberta to encode and incorporate the textual feature is more effective than applying MLP over categorical features.

Table 7: The statistics and hyperparameter settings of the evaluated regular multivariate time series datasets.

| Datasets | Variables | Classes | Length | Train size | Grid layout | Image size | Learning rate | Epochs |
|----------|-----------|---------|--------|-----------|-------------|------------|---------------|--------|
| EC | 3 | 4 | 1,751 | 261 | $2 \times 2$ | $256 \times 256$ | 1e-4 | 20 |
| UW | 3 | 8 | 315 | 120 | $2 \times 2$ | $256 \times 256$ | 1e-4 | 100 |
| SCP1 | 6 | 2 | 896 | 268 | $2 \times 3$ | $256 \times 384$ | 1e-4 | 100 |
| SCP2 | 7 | 2 | 1,152 | 200 | $3 \times 3$ | $384 \times 384$ | 5e-5 | 100 |
| JV | 12 | 9 | 29 | 270 | $4 \times 4$ | $384 \times 384$ | 1e-4 | 100 |
| SAD | 13 | 10 | 93 | 6599 | $4 \times 4$ | $384 \times 384$ | 1e-5 | 20 |
| HB | 61 | 2 | 405 | 204 | $4 \times 4$ | $384 \times 384$ | 1e-4 | 100 |
| FD | 144 | 2 | 62 | 5890 | $12 \times 12$ | $384 \times 384$ | 5e-4 | 100 |
| PS | **963** | 7 | 144 | 267 | $32 \times 32$ | $384 \times 384$ | 5e-4 | 100 |
| EW | 6 | 5 | **17984** | 128 | $2 \times 3$ | $256 \times 384$ | 2e-5 | 100 |

## B.3 Experiment on Regular Time Series

We selected ten representative multivariate time series datasets from the UEA Time Series Classification Archive [1] with diverse characteristics, including the number of classes, variables, and time series length. The datasets we chose are EthanolConcentration (EC), Handwriting (HW), UWaveGestureLibrary (UW), SelfRegulationSCP1 (SCP1), SelfRegulationSCP2 (SCP2), JapaneseVowels (JV), SpokenArabicDigits (SAD), Heartbeat (HB), FaceDetection (FD), PEMS-SF (PS), and EigenWorms

(EW). Notably, the PS dataset has an exceptionally high number of variables (963), while the EW dataset has extremely long time series (17984). These two datasets allow us to assess the effectiveness of our approach when dealing with large numbers of variables and long time series. We applied different image sizes according to the grid layouts for these datasets. The hyperparameter settings are provided in Table 7, and we applied cutout data augmentation methods to SCP1, SCP2, and JV datasets due to the small size of their training sets.

## B.4 Self-supervised Learning

We preliminary explored masked image modeling self-supervised pre-training on the time series line graph images. We randomly mask columns of patches with a width of 32 on each line graph within a grid cell. The masking ratio is set as 50%. We finetuned the Swin Transformer model for 10 epochs with batch size 48. The learning rate is 2e-5. Following [6], we use a linear layer to reconstruct the pixel values and employ an $\ell_1$ loss on the masked pixels:

$$\mathcal{L} = \frac{1}{\Omega(\mathbf{p}_M)} \|\hat{\mathbf{p}_M} - \mathbf{p}_M\|_1 , \tag{1}$$

where $\mathbf{p}_M$ and $\hat{\mathbf{p}_M}$ are the masked and reconstructed pixels, respectively; $\Omega(\cdot)$ denotes the number of elements.

## B.5 Full Experimental Results

We presented the full experimental results in the leave-sensors-out settings in Table 8, and the full results of ablation studies on backbone vision models are presented in Table 9.

Table 8: Full results in the leave-sensors-out settings on PAM dataset. The "missing ratio" denotes the ratio of masked variables.

| Missing ratio | Methods | PAM (Leave-**fixed**-sensors-out) | | | | PAM (Leave-**random**-sensors-out) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Accuracy | Precision | Recall | F1 score | Accuracy | Precision | Recall | F1 score |
| 10% | Transformer | $60.3 \pm 2.4$ | $57.8 \pm 9.3$ | $59.8 \pm 5.4$ | $57.2 \pm 8.0$ | $60.9 \pm 12.8$ | $58.4 \pm 18.4$ | $59.1 \pm 16.2$ | $56.9 \pm 18.9$ |
| | Trans-mean | $60.4 \pm 11.2$ | $61.8 \pm 14.9$ | $60.2 \pm 13.8$ | $58.0 \pm 15.2$ | $62.4 \pm 3.5$ | $59.6 \pm 7.2$ | $63.7 \pm 8.1$ | $62.7 \pm 6.4$ |
| | GRU-D | $65.4 \pm 1.7$ | $72.6 \pm 2.6$ | $64.3 \pm 5.3$ | $63.6 \pm 0.4$ | $68.4 \pm 3.7$ | $74.2 \pm 3.0$ | $70.8 \pm 4.2$ | $72.0 \pm 3.7$ |
| | SeFT | $58.9 \pm 2.3$ | $62.5 \pm 1.8$ | $59.6 \pm 2.6$ | $59.6 \pm 2.6$ | $40.0 \pm 1.9$ | $40.8 \pm 3.2$ | $41.0 \pm 0.7$ | $39.9 \pm 1.5$ |
| | mTAND | $58.8 \pm 2.7$ | $59.5 \pm 5.3$ | $64.4 \pm 2.9$ | $61.8 \pm 4.1$ | $53.4 \pm 2.0$ | $54.8 \pm 2.7$ | $57.0 \pm 1.9$ | $55.9 \pm 2.2$ |
| | Raindrop | $77.2 \pm 2.1$ | $82.3 \pm 1.1$ | $78.4 \pm 1.9$ | $75.2 \pm 3.1$ | $76.7 \pm 1.8$ | $79.9 \pm 1.7$ | $77.9 \pm 2.3$ | $78.6 \pm 1.8$ |
| | **ViTST** | $\mathbf{92.7} \pm 0.9$ | $\mathbf{94.2} \pm 0.9$ | $\mathbf{93.2} \pm 0.4$ | $\mathbf{93.6} \pm 0.6$ | $\mathbf{88.4} \pm 1.4$ | $\mathbf{92.3} \pm 0.5$ | $\mathbf{88.6} \pm 1.9$ | $\mathbf{89.8} \pm 1.5$ |
| 20% | Transformer | $63.1 \pm 7.6$ | $71.1 \pm 7.1$ | $62.2 \pm 8.2$ | $63.2 \pm 8.7$ | $62.3 \pm 11.5$ | $65.9 \pm 12.7$ | $61.4 \pm 13.9$ | $61.8 \pm 15.6$ |
| | Trans-mean | $61.2 \pm 3.0$ | $74.2 \pm 1.8$ | $63.5 \pm 4.4$ | $64.1 \pm 4.1$ | $56.8 \pm 4.1$ | $59.4 \pm 3.4$ | $53.2 \pm 3.9$ | $55.3 \pm 3.5$ |
| | GRU-D | $64.6 \pm 1.8$ | $73.3 \pm 3.6$ | $63.5 \pm 4.6$ | $64.8 \pm 3.6$ | $64.8 \pm 0.4$ | $69.8 \pm 0.8$ | $65.8 \pm 0.5$ | $67.2 \pm 0.0$ |
| | SeFT | $35.7 \pm 0.5$ | $42.1 \pm 4.8$ | $38.1 \pm 1.3$ | $35.0 \pm 2.2$ | $34.2 \pm 2.8$ | $34.9 \pm 5.2$ | $34.6 \pm 2.1$ | $33.3 \pm 2.7$ |
| | mTAND | $33.2 \pm 5.0$ | $36.9 \pm 3.7$ | $37.7 \pm 3.7$ | $37.3 \pm 3.4$ | $45.6 \pm 1.6$ | $49.2 \pm 2.1$ | $49.0 \pm 1.6$ | $49.0 \pm 1.0$ |
| | Raindrop | $66.5 \pm 4.0$ | $72.0 \pm 3.9$ | $67.9 \pm 5.8$ | $65.1 \pm 7.0$ | $71.3 \pm 2.5$ | $75.8 \pm 2.2$ | $72.5 \pm 2.0$ | $73.4 \pm 2.1$ |
| | **ViTST** | $\mathbf{88.4} \pm 1.0$ | $\mathbf{90.4} \pm 1.4$ | $\mathbf{89.3} \pm 0.8$ | $\mathbf{89.7} \pm 1.0$ | $\mathbf{85.1} \pm 1.2$ | $\mathbf{91.1} \pm 1.0$ | $\mathbf{85.6} \pm 1.0$ | $\mathbf{87.0} \pm 1.0$ |
| 30% | Transformer | $31.6 \pm 10.0$ | $26.4 \pm 9.7$ | $24.0 \pm 10.0$ | $19.0 \pm 12.8$ | $52.0 \pm 11.9$ | $55.2 \pm 15.3$ | $50.1 \pm 13.3$ | $48.4 \pm 18.2$ |
| | Trans-mean | $42.5 \pm 8.6$ | $45.3 \pm 9.6$ | $37.0 \pm 7.9$ | $33.9 \pm 8.2$ | $65.1 \pm 1.9$ | $63.8 \pm 1.2$ | $67.9 \pm 1.8$ | $64.9 \pm 1.7$ |
| | GRU-D | $45.1 \pm 2.9$ | $51.7 \pm 6.2$ | $42.1 \pm 6.6$ | $47.2 \pm 3.9$ | $58.0 \pm 2.0$ | $63.2 \pm 1.7$ | $58.2 \pm 3.1$ | $55.3 \pm 3.5$ |
| | SeFT | $32.7 \pm 2.3$ | $27.9 \pm 2.4$ | $34.5 \pm 3.0$ | $28.0 \pm 1.4$ | $31.7 \pm 1.5$ | $31.0 \pm 2.7$ | $32.0 \pm 1.2$ | $28.0 \pm 1.6$ |
| | mTAND | $27.5 \pm 4.5$ | $31.2 \pm 7.3$ | $30.6 \pm 4.0$ | $30.8 \pm 5.6$ | $34.7 \pm 5.5$ | $43.4 \pm 4.0$ | $36.3 \pm 4.7$ | $39.5 \pm 4.4$ |
| | Raindrop | $52.4 \pm 2.8$ | $60.9 \pm 3.8$ | $51.3 \pm 7.1$ | $48.4 \pm 1.8$ | $60.3 \pm 3.5$ | $68.1 \pm 3.1$ | $60.3 \pm 3.6$ | $61.9 \pm 3.9$ |
| | **ViTST** | $\mathbf{84.1} \pm 1.3$ | $\mathbf{86.5} \pm 0.4$ | $\mathbf{83.1} \pm 0.8$ | $\mathbf{84.9} \pm 1.0$ | $\mathbf{80.6} \pm 1.2$ | $\mathbf{89.5} \pm 1.3$ | $\mathbf{80.9} \pm 1.1$ | $\mathbf{82.6} \pm 1.1$ |
| 40% | Transformer | $23.0 \pm 3.5$ | $7.4 \pm 6.0$ | $14.5 \pm 2.6$ | $6.9 \pm 2.6$ | $43.8 \pm 14.0$ | $44.6 \pm 23.0$ | $40.5 \pm 15.9$ | $40.2 \pm 20.1$ |
| | Trans-mean | $25.7 \pm 2.5$ | $9.1 \pm 2.3$ | $18.5 \pm 1.4$ | $9.9 \pm 1.1$ | $48.7 \pm 2.7$ | $55.8 \pm 2.6$ | $54.2 \pm 3.0$ | $55.1 \pm 2.9$ |
| | GRU-D | $46.4 \pm 2.5$ | $64.5 \pm 6.8$ | $42.6 \pm 7.4$ | $44.3 \pm 7.9$ | $47.7 \pm 1.4$ | $63.4 \pm 1.6$ | $44.5 \pm 0.5$ | $47.5 \pm 0.0$ |
| | SeFT | $26.3 \pm 0.9$ | $29.9 \pm 4.5$ | $27.3 \pm 1.6$ | $22.3 \pm 1.9$ | $26.8 \pm 2.6$ | $2.41 \pm 3.4$ | $28.0 \pm 1.2$ | $23.3 \pm 3.0$ |
| | mTAND | $19.4 \pm 4.5$ | $15.1 \pm 4.4$ | $20.2 \pm 3.8$ | $17.0 \pm 3.4$ | $23.7 \pm 1.0$ | $33.9 \pm 6.5$ | $26.4 \pm 1.6$ | $29.3 \pm 1.9$ |
| | Raindrop | $52.5 \pm 3.7$ | $53.4 \pm 5.6$ | $48.6 \pm 1.9$ | $44.7 \pm 3.4$ | $57.0 \pm 3.1$ | $65.4 \pm 2.7$ | $56.7 \pm 3.1$ | $58.9 \pm 2.5$ |
| | **ViTST** | $\mathbf{76.5} \pm 1.9$ | $\mathbf{83.5} \pm 0.9$ | $\mathbf{76.7} \pm 2.4$ | $\mathbf{78.3} \pm 2.1$ | $\mathbf{73.7} \pm 2.2$ | $\mathbf{86.4} \pm 1.1$ | $\mathbf{74.0} \pm 2.2$ | $\mathbf{75.8} \pm 1.8$ |
| 50% | Transformer | $21.4 \pm 1.8$ | $2.7 \pm 0.2$ | $12.5 \pm 0.4$ | $4.4 \pm 0.3$ | $43.2 \pm 2.5$ | $52.0 \pm 2.5$ | $36.9 \pm 3.1$ | $41.9 \pm 3.2$ |
| | Trans-mean | $21.3 \pm 1.6$ | $2.8 \pm 0.4$ | $12.5 \pm 0.7$ | $4.6 \pm 0.2$ | $46.4 \pm 1.4$ | $59.1 \pm 3.2$ | $43.1 \pm 2.2$ | $46.5 \pm 3.1$ |
| | GRU-D | $37.3 \pm 2.7$ | $29.6 \pm 5.9$ | $32.8 \pm 4.6$ | $26.6 \pm 5.9$ | $49.7 \pm 1.2$ | $52.4 \pm 0.3$ | $42.5 \pm 1.7$ | $47.5 \pm 1.2$ |
| | SeFT | $24.7 \pm 1.7$ | $15.9 \pm 2.7$ | $25.3 \pm 2.6$ | $18.2 \pm 2.4$ | $26.4 \pm 1.4$ | $23.0 \pm 2.9$ | $27.5 \pm 0.4$ | $23.5 \pm 1.8$ |
| | mTAND | $16.9 \pm 3.1$ | $12.6 \pm 5.5$ | $17.0 \pm 1.6$ | $13.9 \pm 4.0$ | $20.9 \pm 3.1$ | $35.1 \pm 6.1$ | $23.0 \pm 3.2$ | $27.7 \pm 3.9$ |
| | Raindrop | $46.6 \pm 2.6$ | $44.5 \pm 2.6$ | $42.4 \pm 3.9$ | $38.0 \pm 4.0$ | $47.2 \pm 4.4$ | $59.4 \pm 3.9$ | $44.8 \pm 5.3$ | $47.6 \pm 5.2$ |
| | **ViTST** | $\mathbf{70.0} \pm 2.7$ | $\mathbf{79.9} \pm 2.2$ | $\mathbf{70.5} \pm 3.1$ | $\mathbf{72.2} \pm 3.0$ | $\mathbf{70.9} \pm 1.2$ | $\mathbf{83.6} \pm 2.4$ | $\mathbf{71.5} \pm 1.4$ | $\mathbf{73.3} \pm 2.1$ |

Table 9: Full results of our approach with different backbone vision models and the compared baselines. **Bold** indicates the best performer, while <u>underline</u> represents the second best.

| Methods | P19 | | P12 | | PAM | | | |
|---|---|---|---|---|---|---|---|---|
| | AUROC | AUPRC | AUROC | AUPRC | Accuracy | Precision | Recall | F1 score |
| Transformer | $80.7 \pm 3.8$ | $42.7 \pm 7.7$ | $83.3 \pm 0.7$ | $47.9 \pm 3.6$ | $83.5 \pm 1.5$ | $84.8 \pm 1.5$ | $86.0 \pm 1.2$ | $85.0 \pm 1.3$ |
| Trans-mean | $83.7 \pm 1.8$ | $45.8 \pm 3.2$ | $82.6 \pm 2.0$ | $46.3 \pm 4.0$ | $83.7 \pm 2.3$ | $84.9 \pm 2.6$ | $86.4 \pm 2.1$ | $85.1 \pm 2.4$ |
| GRU-D | $83.9 \pm 1.7$ | $46.9 \pm 2.1$ | $81.9 \pm 2.1$ | $46.1 \pm 4.7$ | $83.3 \pm 1.6$ | $84.6 \pm 1.2$ | $85.2 \pm 1.6$ | $84.8 \pm 1.2$ |
| SeFT | $81.2 \pm 2.3$ | $41.9 \pm 3.1$ | $73.9 \pm 2.5$ | $31.1 \pm 4.1$ | $67.1 \pm 2.2$ | $70.0 \pm 2.4$ | $68.2 \pm 1.5$ | $68.5 \pm 1.8$ |
| mTAND | $84.4 \pm 1.3$ | $50.6 \pm 2.0$ | $84.2 \pm 0.8$ | $48.2 \pm 3.4$ | $74.6 \pm 4.3$ | $74.3 \pm 4.0$ | $79.5 \pm 2.8$ | $76.8 \pm 3.4$ |
| IP-Net | $84.6 \pm 1.3$ | $38.1 \pm 3.7$ | $82.6 \pm 1.4$ | $47.6 \pm 3.1$ | $74.3 \pm 3.8$ | $75.6 \pm 2.1$ | $77.9 \pm 2.2$ | $76.6 \pm 2.8$ |
| DGM$^2$-O | $86.7 \pm 3.4$ | $44.7 \pm 11.7$ | $84.4 \pm 1.6$ | $47.3 \pm 3.6$ | $82.4 \pm 2.3$ | $85.2 \pm 1.2$ | $83.9 \pm 2.3$ | $84.3 \pm 1.8$ |
| MTGNN | $81.9 \pm 6.2$ | $39.9 \pm 8.9$ | $74.4 \pm 6.7$ | $35.5 \pm 6.0$ | $83.4 \pm 1.9$ | $85.2 \pm 1.7$ | $86.1 \pm 1.9$ | $85.9 \pm 2.4$ |
| Raindrop | $87.0 \pm 2.3$ | <u>$51.8 \pm 5.5$</u> | $82.8 \pm 1.7$ | $44.0 \pm 3.0$ | $88.5 \pm 1.5$ | $89.9 \pm 1.5$ | $89.9 \pm 0.6$ | $89.8 \pm 1.0$ |
| ResNet | $76.3 \pm 3.3$ | $34.7 \pm 4.1$ | $72.9 \pm 1.0$ | $28.8 \pm 2.4$ | $73.1 \pm 0.9$ | $82.4 \pm 5.6$ | $69.7 \pm 0.9$ | $71.4 \pm 1.8$ |
| ViT | <u>$87.9 \pm 2.5$</u> | $51.6 \pm 3.7$ | <u>$84.8 \pm 1.3$</u> | $48.1 \pm 3.8$ | <u>$93.4 \pm 0.7$</u> | <u>$94.7 \pm 0.9$</u> | <u>$94.1 \pm 0.7$</u> | <u>$94.3 \pm 0.7$</u> |
| **Swin** | **$89.4 \pm 1.9$** | **$52.8 \pm 3.8$** | **$85.6 \pm 1.1$** | $49.8 \pm 2.5$ | **$96.1 \pm 0.7$** | **$96.8 \pm 1.1$** | **$96.5 \pm 0.7$** | **$96.6 \pm 0.9$** |
| Swin-scratch | $74.6 \pm 2.5$ | $29.9 \pm 4.6$ | $66.9 \pm 1.6$ | $26.5 \pm 2.6$ | $84.5 \pm 0.5$ | $86.6 \pm 0.6$ | $87.1 \pm 1.2$ | $86.6 \pm 0.6$ |

# References

[1] Bagnall, A., Dau, H. A., Lines, J., Flynn, M., Large, J., Bostrom, A., Southam, P., and Keogh, E. The uea multivariate time series classification archive, 2018. *arXiv preprint arXiv:1811.00075*, 2018.

[2] Goldberger, A. L., Amaral, L. A., Glass, L., Hausdorff, J. M., Ivanov, P. C., Mark, R. G., Mietus, J. E., Moody, G. B., Peng, C.-K., and Stanley, H. E. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *circulation*, 101(23): e215–e220, 2000.

[3] OpenAI. Gpt-4 technical report, 2023.

[4] Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pp. 8748–8763. PMLR, 2021.

[5] Reyna, M. A., Josef, C., Seyedi, S., Jeter, R., Shashikumar, S. P., Westover, M. B., Sharma, A., Nemati, S., and Clifford, G. D. Early prediction of sepsis from clinical data: the physionet/computing in cardiology challenge 2019. In *2019 Computing in Cardiology (CinC)*, pp. Page–1. IEEE, 2019.

[6] Xie, Z., Zhang, Z., Cao, Y., Lin, Y., Bao, J., Yao, Z., Dai, Q., and Hu, H. Simmim: A simple framework for masked image modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9653–9663, 2022.

[7] Zhang, X., Zeman, M., Tsiligkaridis, T., and Zitnik, M. Graph-guided network for irregularly sampled multivariate time series. In *International Conference on Learning Representations*, 2022.