# Weak-Mamba-UNet

Marijan Sorić
Žan Stanonik

NTNU – Computer Vision and Deep Learning – TDT4265

SOTA Paper Challenge

# Article presentation

- ○ Medical Image Segmentation
- ○ Preprint: 16 Feb 2024

---

## Weak-Mamba-UNet:
## Visual Mamba Makes CNN and ViT Work Better
## for Scribble-based Medical Image Segmentation
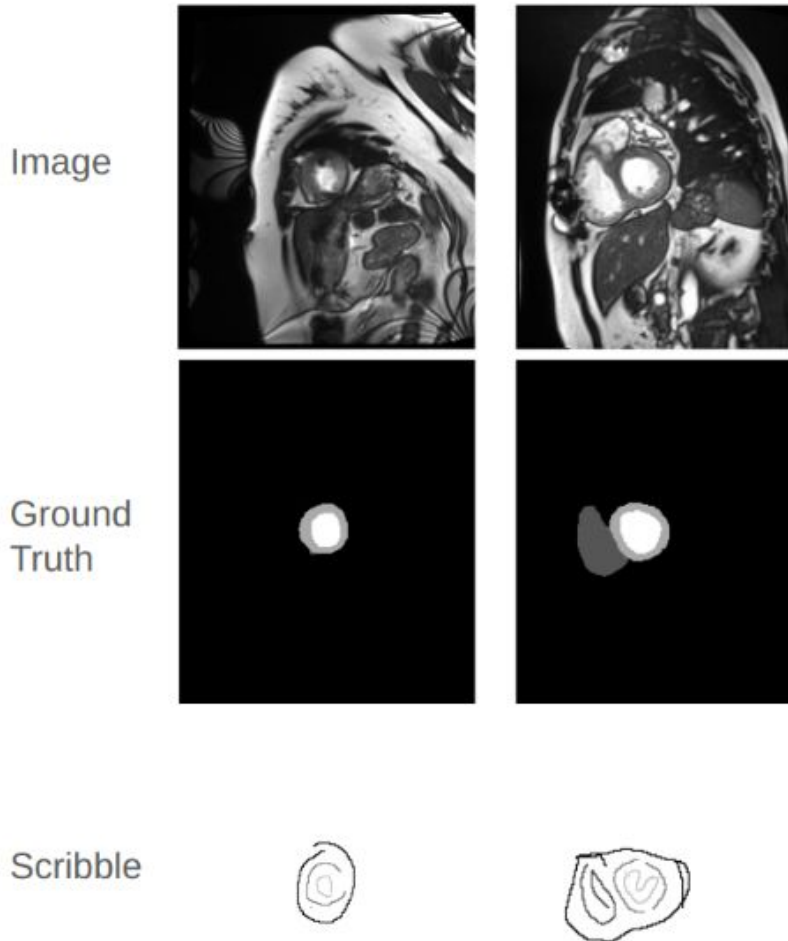
Ziyang Wang[1] and Chao Ma[2]

[1] Department of Computer Science, University of Oxford, UK
[2] Mianyang Visual Object Detection and Recognition Engineering Center, China
ziyang.wang@cs.ox.ac.uk
https://github.com/ziyangwang007/Mamba-UNet

# Introduction



- Method: CNN, ViT & Visual Mamba

- Combines deep learning with the efficiency of WSL

- Architectures benefit from each other

# Results

## Benchmarks

- Outperforms single architectures

- Showcases advantages of:
  - limited supervision
  - limited resources

**Table 2.** Ablation Studies on Different Combinations of Segmentation Backbone Networks with the Same WSL Framework.

| Network | Dice↑ | Acc↑ | Pre↑ | Sen↑ | Spe↑ | HD↓ | ASD↓ |
|---|---|---|---|---|---|---|---|
| 3×UNet | 0.9141 | 0.9959 | 0.8958 | 0.9383 | 0.9927 | 8.0566 | 2.8806 |
| 3×SwinUNet | 0.7446 | 0.9791 | 0.6555 | 0.9142 | 0.9815 | 121.4224 | 51.4317 |
| 3×MambaUNet | 0.9128 | 0.9958 | 0.8931 | 0.9395 | 0.9932 | 8.3386 | 2.7928 |
| UNet+SwinUNet+MambaUNet(Ours) | 0.9171 | 0.9963 | 0.9095 | 0.9309 | 0.9920 | 3.9597 | 0.8810 |

# The main idea

- Three distinct architectures

$$\mathbf{Y}_{\text{pseudo}} = \alpha \times f_{\text{cnn}}(\mathbf{X};\theta) + \beta \times f_{\text{vit}}(\mathbf{X};\theta) + \gamma \times f_{\text{mamba}}(\mathbf{X};\theta)$$

$$\alpha + \beta + \gamma = 1$$

- Multi-view cross-SL

$$\mathcal{L}_{\text{total}} = \sum_{i=1}^{3}(\mathcal{L}_{\text{pce}}^{i} + \mathcal{L}_{\text{dice}}^{i})$$

- Overall loss
  - Scribble-based

$$\mathcal{L}_{\text{pce}} = -\sum_{i \in \Omega_L}\sum_{k} y_{\text{s}}[i,k]\log(y_{\text{p}}[i,k])$$ (Pixels annotated with scribbles)

  - Dense-signal pseudo label

$$\mathcal{L}_{\text{dice}} = \text{Dice}\big(\text{argmax}(f(\boldsymbol{X};\theta), \boldsymbol{Y}_{\text{pseudo}})\big)$$

# Architecture
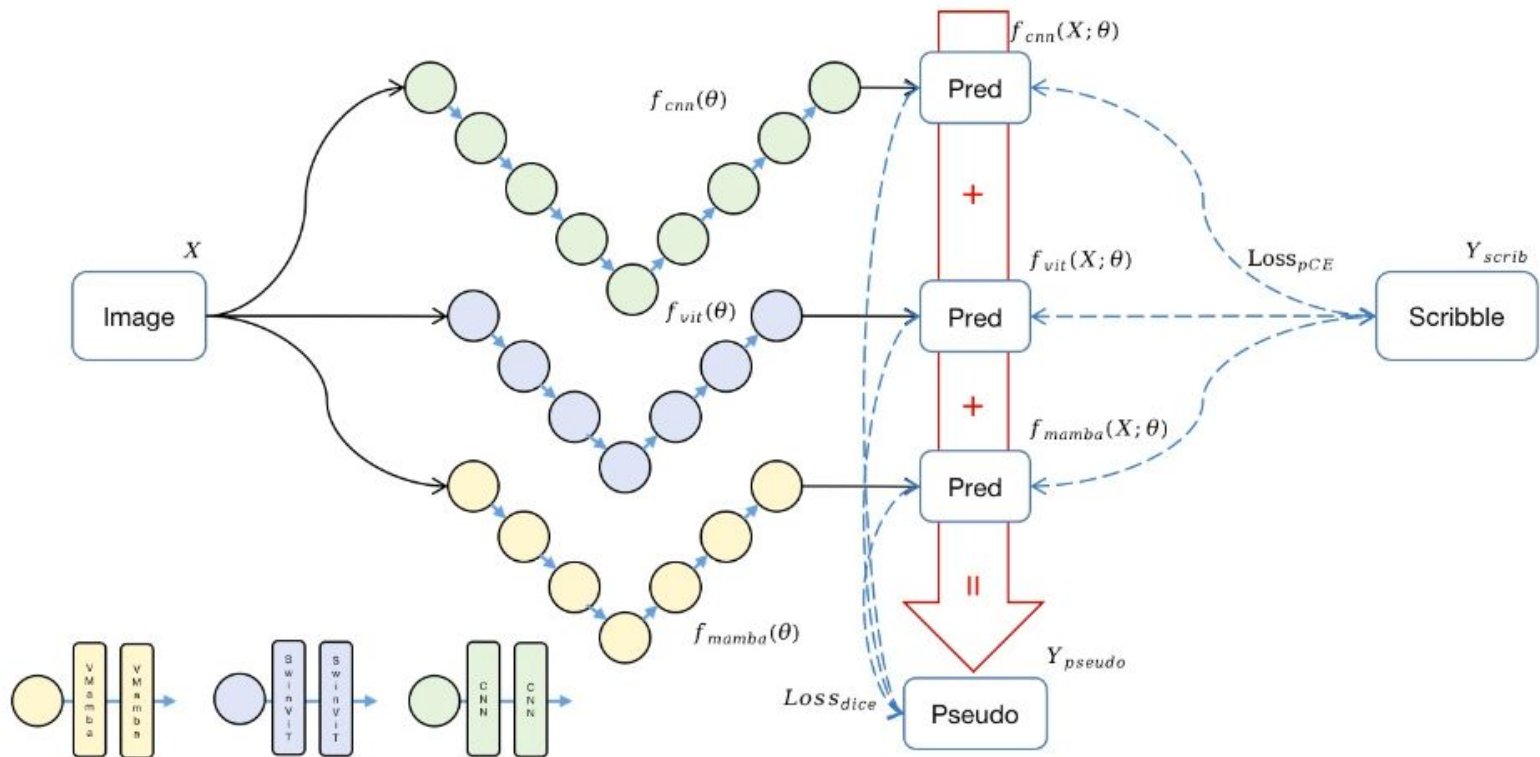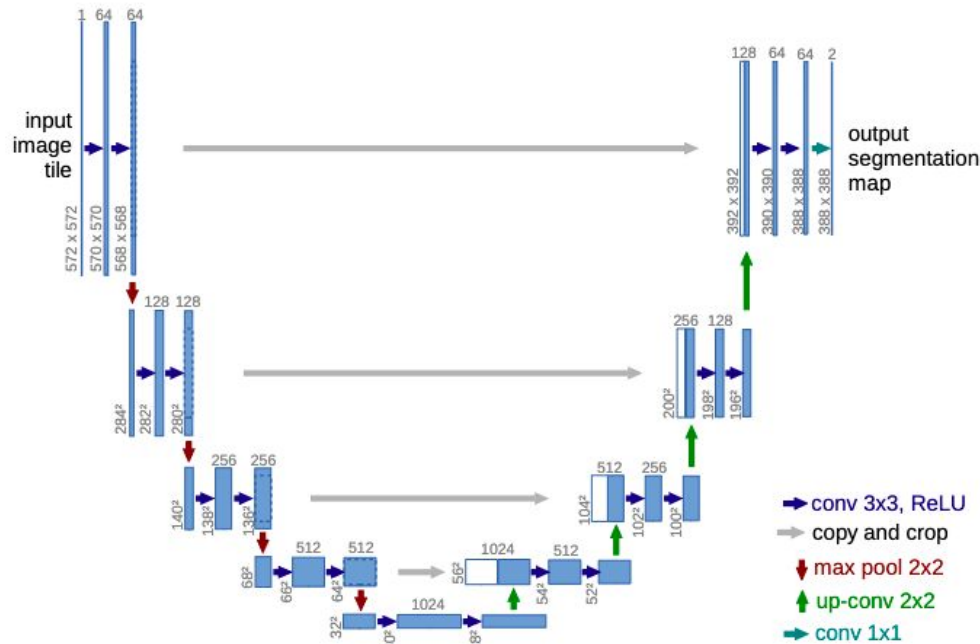
## Overview of the model



**Fig. 2.** Semi-Mamba-UNet: The Framework of Contrastive Cross-Supervised Visual Mamba-based UNet for Semi-Supervised Medical Image Segmentation.
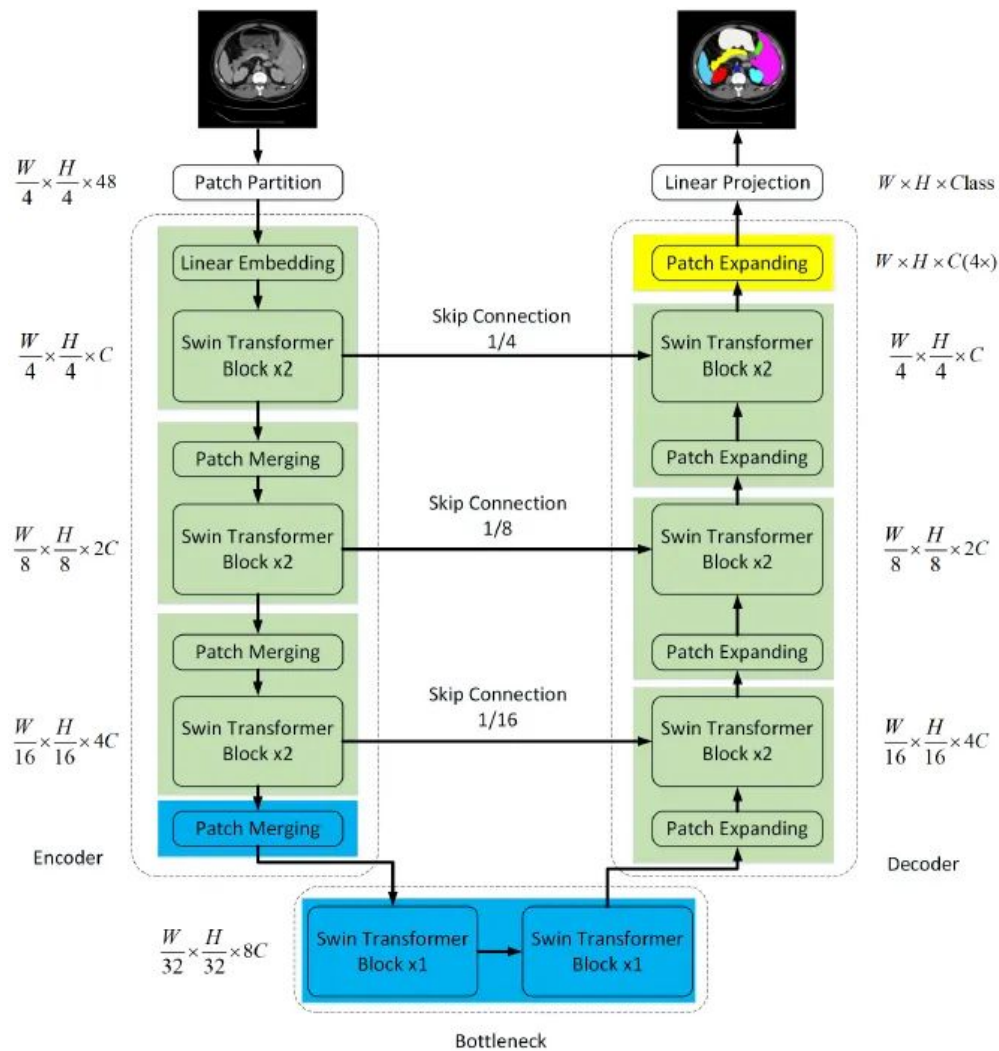
# CNN-based UNet



Fig. 1. U-net architecture (example for 32x32 pixels in the lowest resolution). Each blue box corresponds to a multi-channel feature map. The number of channels is denoted on top of the box. The x-y-size is provided at the lower left edge of the box. White boxes represent copied feature maps. The arrows denote the different operations.

- Encoder: What?

- Decoder: Where?

- Skip connection

# SwinUNet

- Variant of UNet

- Swin captures long-range dependencies

- Image is broken down to patches

- UNet decoder-encoder condenses information
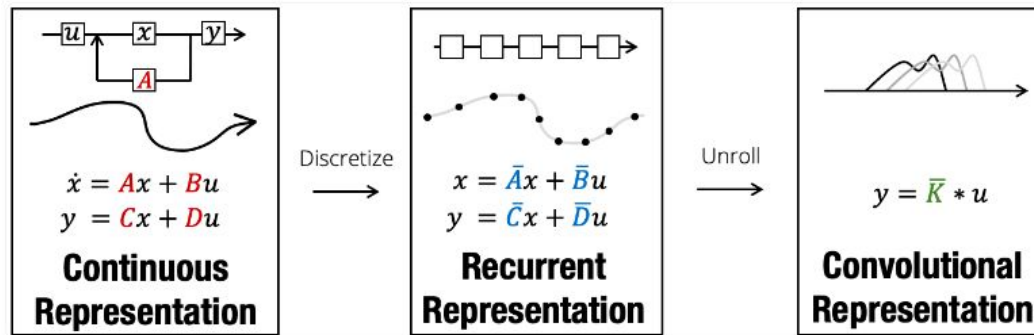
- Ideal for medical image segmentation

# SwinUNet

# Mamba explained

## RNN → SSM → S4 →S6 *(Mamba)*

- State Space Models

- Linear computation



$$h'(t) = Ah(t) + Bx(t) \quad (1a) \qquad h_t = \overline{A}h_{t-1} + \overline{B}x_t \quad (2a) \qquad \overline{K} = (C\overline{B}, C\overline{AB}, \dots, C\overline{A}^k\overline{B}, \dots) \quad (3a)$$

$$y(t) = Ch(t) \quad (1b) \qquad y_t = Ch_t \quad (2b) \qquad y = x * \overline{K} \quad (3b)$$

| **Algorithm 1** SSM (S4) | **Algorithm 2** SSM + Selection (S6) |
|---|---|
| **Input:** $x : (B, L, D)$ | **Input:** $x : (B, L, D)$ |
| **Output:** $y : (B, L, D)$ | **Output:** $y : (B, L, D)$ |
| 1: $A : (D, N) \leftarrow$ Parameter | 1: $A : (D, N) \leftarrow$ Parameter |
| $\triangleright$ Represents structured $N \times N$ matrix | $\triangleright$ Represents structured $N \times N$ matrix |
| 2: $B : (D, N) \leftarrow$ Parameter | 2: $B : (B, L, N) \leftarrow s_B(x)$ |
| 3: $C : (D, N) \leftarrow$ Parameter | 3: $C : (B, L, N) \leftarrow s_C(x)$ |
| 4: $\Delta : (D) \leftarrow \tau_\Delta$(Parameter) | 4: $\Delta : (B, L, D) \leftarrow \tau_\Delta$(Parameter$+s_\Delta(x)$) |
| 5: $\overline{A}, \overline{B} : (D, N) \leftarrow$ discretize$(\Delta, A, B)$ | 5: $\overline{A}, \overline{B} : (B, L, D, N) \leftarrow$ discretize$(\Delta, A, B)$ |
| 6: $y \leftarrow$ SSM$(\overline{A}, \overline{B}, C)(x)$ | 6: $y \leftarrow$ SSM$(\overline{A}, \overline{B}, C)(x)$ |
| $\triangleright$ Time-invariant: recurrence or convolution | $\triangleright$ Time-varying: recurrence (*scan*) only |
| 7: **return** $y$ | 7: **return** $y$ |

# Mamba-based MambaUNet



Fig. 2. The architecture of Mamba-UNet, which is composed of encoder, bottleneck, decoder and skip connections. The encoder, bottleneck and decoder are all constructed based on Visual Mamba block.

- Encoder-Decoder architecture

- Efficient long-range dependency modelling



VSS Block

# Conclusion

- Reduces the cost & resources required for annotation

- Different algorithms completed each other

- Simpler form of annotation

- Can be applied to a variety of ML analysis tasks

# Sources

## Papers:

- Weak-Mamba-UNet: Visual Mamba Makes CNN and ViT Work Better for Scribble-based Medical Image Segmentation
- U-Net: Convolutional Networks for Biomedical Image Segmentation
- Swin-Unet: Unet-like Pure Transformer for Medical Image Segmentation
- Mamba: Linear-Time Sequence Modeling with Selective State Spaces
- Mamba-UNet: UNet-Like Pure Visual Mamba for Medical Image Segmentation

## Blogs:

- Introduction to State Space Models (SSM) – *Hugging Face*