



# SpikingSSMs: Learning Long Sequences with Sparse and Parallel Spiking State Space Models

Shuaijie Shen<sup>\*,1,2</sup>, Chao Wang<sup>\*,1,2</sup>, Renzhuo Huang<sup>1,2</sup>, Yan Zhong<sup>2,3</sup>, Qinghai Guo<sup>2</sup>, Zhichao Lu<sup>4</sup>,  
Jianguo Zhang<sup>†,1,5</sup>, Luziwei Leng<sup>†,2</sup>

<sup>1</sup> SUSTech | <sup>2</sup> Huawei | <sup>3</sup> Peking University | <sup>4</sup> CityU | <sup>5</sup> Pengcheng Lab



# State Space Models (SSMs)

$$\begin{aligned}h_t &= Ah_{t-1} + Bx_t \\y_t &= Ch_t + Dx_t\end{aligned}$$

**SSMs**

Utilizing FFT for parallel training  
with  $\mathcal{O}(L \log L)$  complexity.

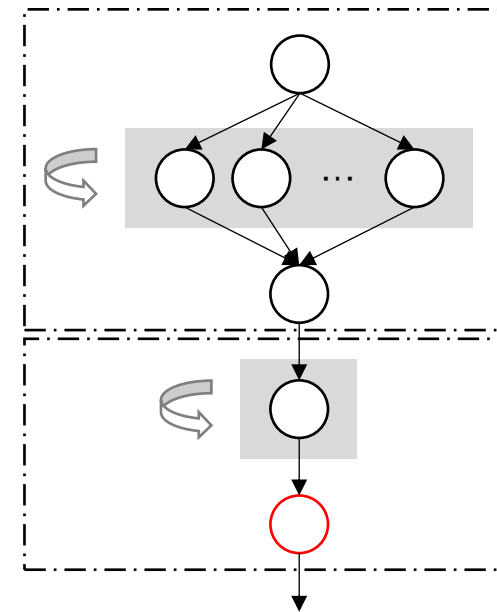
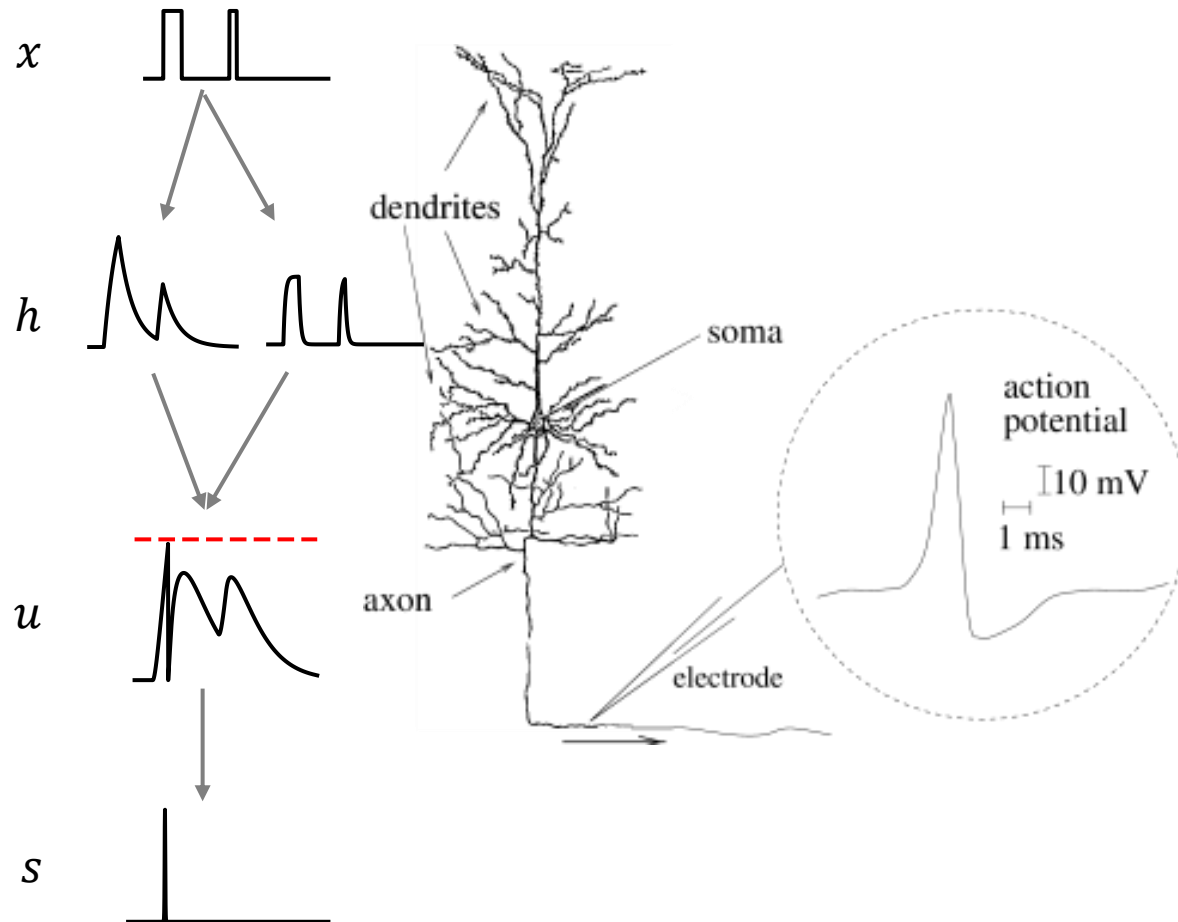
**Parallel via FFT**

e.g., S4、S4D

**Long-Range  
Dependency Modeling**

# SSMs Meet Dendritic Computation

**Biological Inspiration:** Dendritic neurons process multi-timescale inputs.



SSM  $\approx$  Dendrites (integration)

LIF  $\approx$  Soma (sparse output)

**SpikingSSMs Design:**

SSM hidden states ( $h$ ) model dendritic integration.  
LIF neuron acts as soma for sparse spiking.

# Parallel Challenge of Spiking Neurons

$$\begin{aligned} u'_t &= \tau u_{t-1} + y_t \\ s_t &= \begin{cases} 1, & u'_t \geq v_{th} \\ 0, & u'_t < v_{th} \end{cases} \\ u_t &= (1 - s_t)u'_t + s_tv_r \end{aligned} \quad \text{Nonlinear !}$$

**LIF (Leaky Integrate and Fire) neuron**  
with hard reset mechanism

**Problem:** LIF neurons require iterative reset,  
limits parallel training.

**Note:**  $s_{1:T} = f_{\theta}(y_{1:T}; v_{th})$

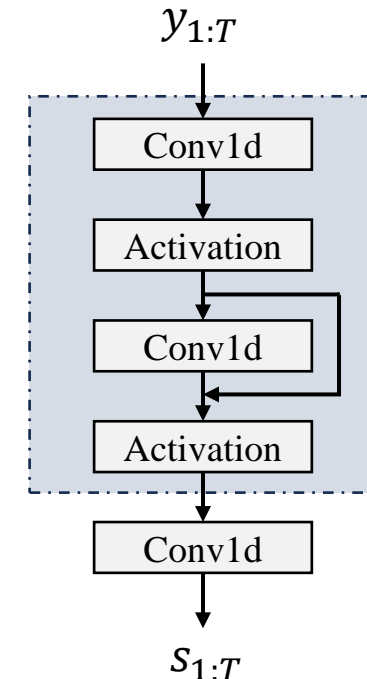
↓  
**A model?**

## Surrogate Dynamic Network (SDN)

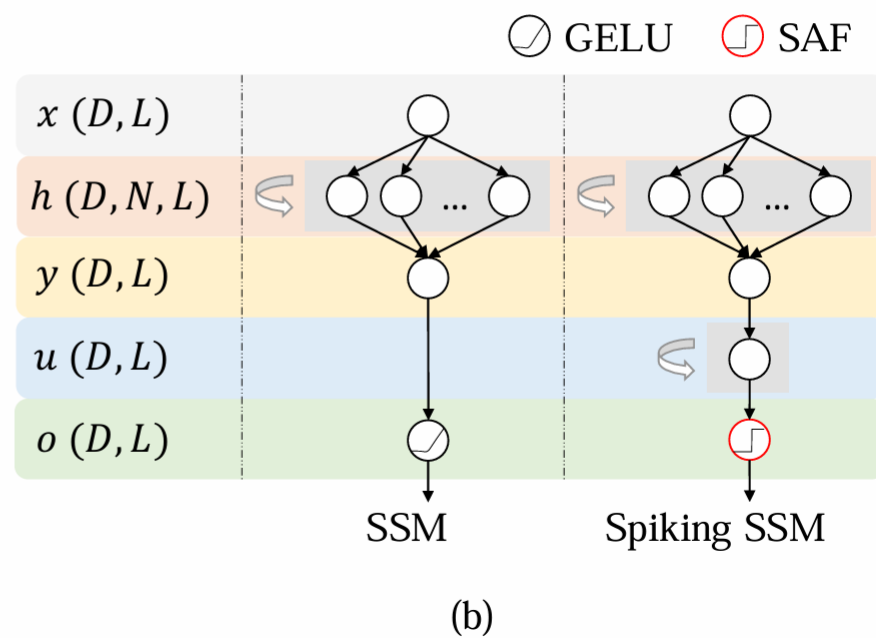
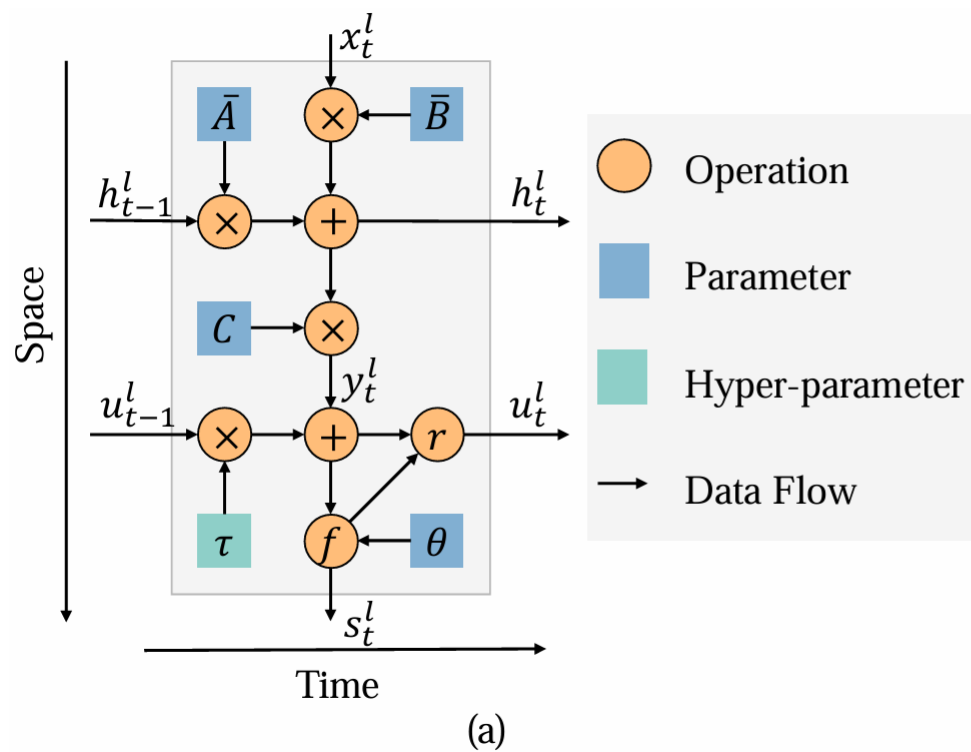
- Lightweight CNN predicts spikes in parallel.
- Architecture: 3-layer 1D-CNN (<200 params).

## Advantage:

- Parallel training
- Can be Removed in iterative inference



# Architecture

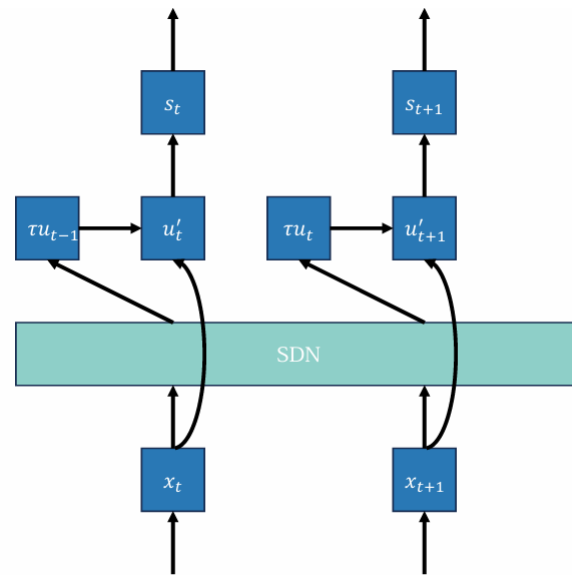


# Computational Graph

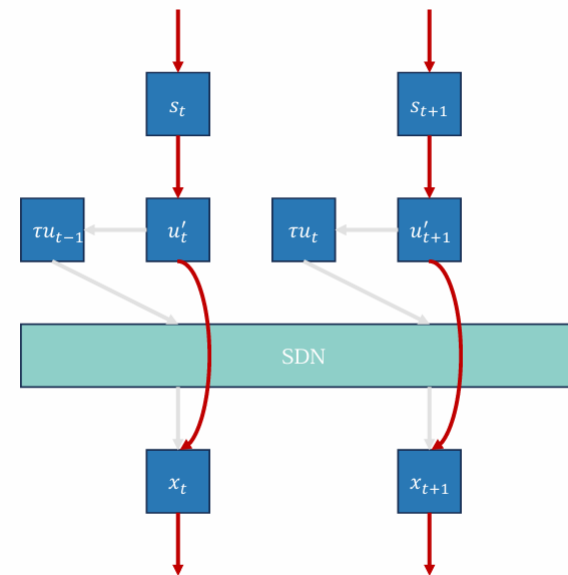
Complex computational graph for SDN in training task models

**Solution:**

$$s_{1:T} = f_{\theta}(y_{1:T}; v_{th}) \longrightarrow \begin{aligned} &\tau u_{0:T-1} = f_{\theta}(y_{1:T}; v_{th}) \quad \text{without gradient propagation} \\ &s_{1:T} = H(\tau u'_{0:T-1} + x_{1:T} - v_{th}) \end{aligned}$$



**Forward of SDN**



**Backward of SDN**

# Learnable Threshold

Optimizing the threshold can improve network performance.

Can SDN approximate neuron dynamics with different threshold?

If both the initial membrane potential and the reset potential are 0:

**Property 1.** The ratio of inputs and threshold determines the dynamic process of the neuron.

$$s_{1:T} = f_{\theta}(y_{1:T}; v_{th}) = f_{\theta}(\alpha y_{1:T}; \alpha v_{th})$$

**Property 2.** The threshold scales the distribution of the input.

$$s_{1:T} = f_{\theta}(y_{1:T}; v_{th}) = f_{\theta}\left(\frac{y_{1:T}}{v_{th}}; 1\right)$$

# SDN Accelerates Training by 100×

Method	Speed			
	L = 1K	L = 2K	L = 4K	L = 8K
BPTT	1370	2900	8040	25600
SLTT	1210	2720	7740	25600
Ours	183	196	200	253
Ratio	7.5 ×	15.0 ×	40.2 ×	101.2 ×

**101× faster than BPTT.**

**Table1: Comparison on training speed of different methods. The input has a batch size of 64. Training with SDN achieves significant acceleration, the speed up ratio amplifies with increasing sequence length.**

Model	Accuracy (%)	Sparsity (%)	Speed (ms)
LIF	85.45	12.08	1480
SDN	85.57	11.92	230
SDN-S	81.52	18.30	285

**SDN matches iterative LIF (85.6% vs. 85.5% accuracy).**

**Table2: Performance comparison on the sCIFAR10 dataset.**

Fixed  
Train from scratch



# Results

	Model	SNN	LISTOPS	TEXT	RETRIEVAL	IMAGE	PATHFINDER	Path-X	AVG
Fixed threshold	Transformer	No	36.37	64.27	57.46	42.44	71.40	—	53.66
	LMUFormer	No	34.43	68.27	78.65	54.16	69.90	—	59.24
	S4D-Lin	No	<b>60.52</b>	<b>86.97</b>	<b>90.96</b>	<b>87.93</b>	<b>93.96</b>	<b>92.80*</b>	<b>85.52</b>
	Spiking LMUFormer	Yes	37.30	65.80	79.76	55.65	72.68	—	60.20
	Binary S4D	Yes	54.80	<b>82.50</b>	85.03	82.00	82.60	61.20	74.69
	S6-based SNN	Yes	55.70	77.62	88.48	80.10	83.41	—	72.55
	SpikingSSM-VF	Yes	59.93	82.35	88.20	86.81	<b>93.68</b>	94.80	84.30
	(spiking rate)		(0.13)	(0.10)	(0.06)	(0.22)	(0.07)	(0.10)	(0.11)
	SpikingSSM-VT	Yes	<b>60.23</b>	80.41	<b>88.77</b>	<b>88.21</b>	93.51	<b>94.82</b>	<b>84.33</b>
	(spiking rate)		(0.14)	(0.06)	(0.06)	(0.15)	(0.08)	(0.10)	(0.10)

Table 3: Performance comparison of SpikingSSM and previous works on the LRA dataset. \*Since the original S4D-Lin failed in the Path-X task, we instead present the result of another close variant S4D-Inv. -VF and -VT denote fixed and trainable threshold, respectively. Furthermore, we take the 50% accuracy for the absence of Path-X accuracy as did in the work of S4D, then compute the overall average metrics across all tasks as AVG. The spiking rate for each task have also been calculated, which is indicated by blue font.

Model	SNN	PPL	Parameters
Transformer	No	20.51	231M
S4	No	20.95	249M
SpikeGPT	Yes	39.75	213M
SpikingSSM	Yes	33.94	75M

Table 4: Performance comparison of SpikingSSMs with previous works on WikiText-103 dataset.



# Energy Efficiency

**Spikes  $\times$  Weights (in SNN):** synaptic Accumulation(AC)

**1  $\times$  float :** + weight

**0  $\times$  float :** no operation

**Activations  $\times$  Weights (in ANN):** Multiply-and-Accumulate(MAC)

**Major Computation:** feature-mix layers (after SSMs)

For Wiki-Text-103 with  $L = 8K$  and 16 layers,  
projecting inputs from  $d = 1024$  to  $d = 2048$   
sparsity: 73.6 %

**Done via multiplication(ANN):**     $275.2G$  MAC     $1.265J$

**Done via accumulation(SNN):**     $72.66G$  AC     $65.40 mJ$

**Theoretical Saving: 95% energy reduction.**

FP	
FAdd	
16 bit	0.4pJ
32 bit	0.9pJ
FMult	
16 bit	1.1pJ
32 bit	3.7pJ

# Conclusion & Future Work

- **Key Takeaways:**
  - SpikingSSMs: Combines SSMs' performance with SNNs' sparsity.
  - SDN: Enables 100× faster training without sacrificing accuracy.
- **Future Work:**
  - Scaling to billion-parameter models.
  - Exploring Brain-Inspired Large Language Models.



**Thank You! Questions?**