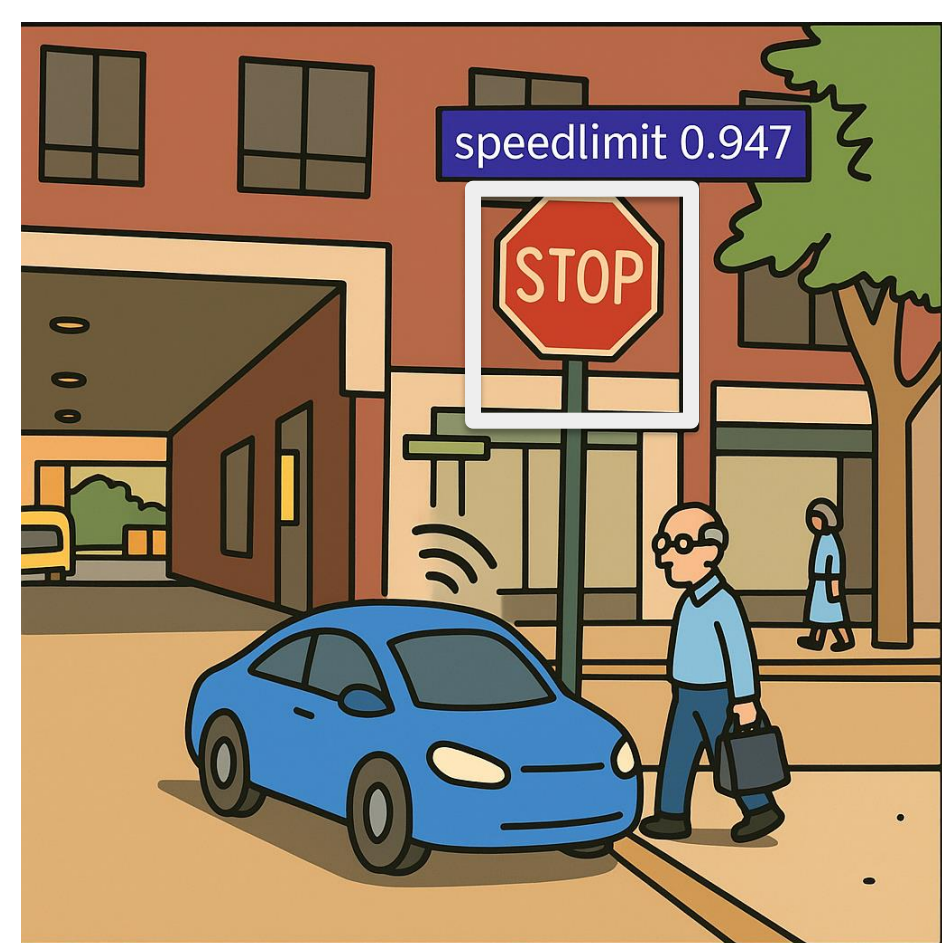


# Unsupervised Backdoor Detection and Mitigation for Spiking Neural Networks

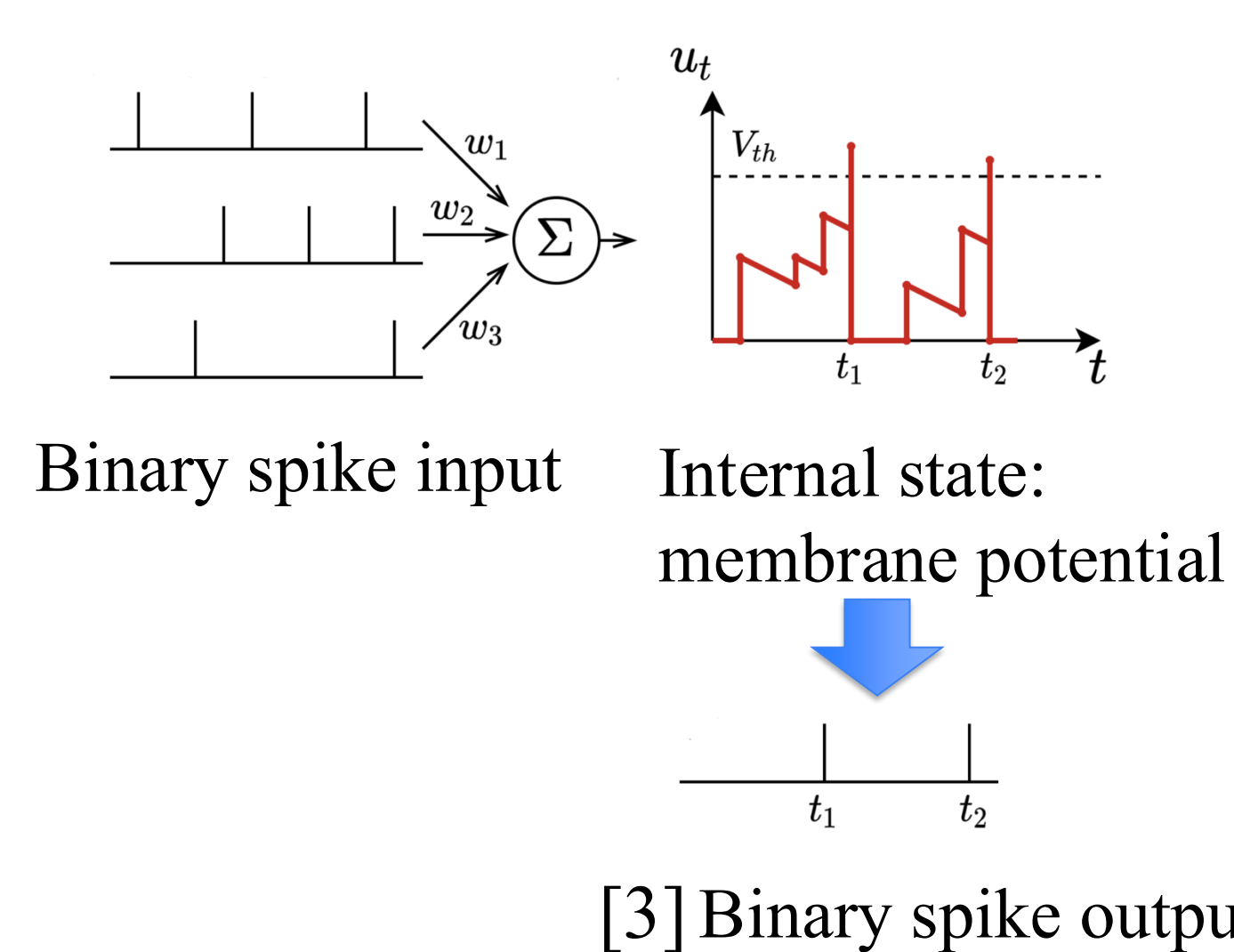
Jiachen Li, Bang Wu, Xiaoyu Xia, Xiaoning Liu, Xun Yi, Xiuzhen Zhang  
RMIT University

## Motivation & Problem

- SNNs used in safety-critical systems (e.g. driving [1]).
- ANN defenses fail due to event-driven, binary spike, threshold nonlinearities of SNN.

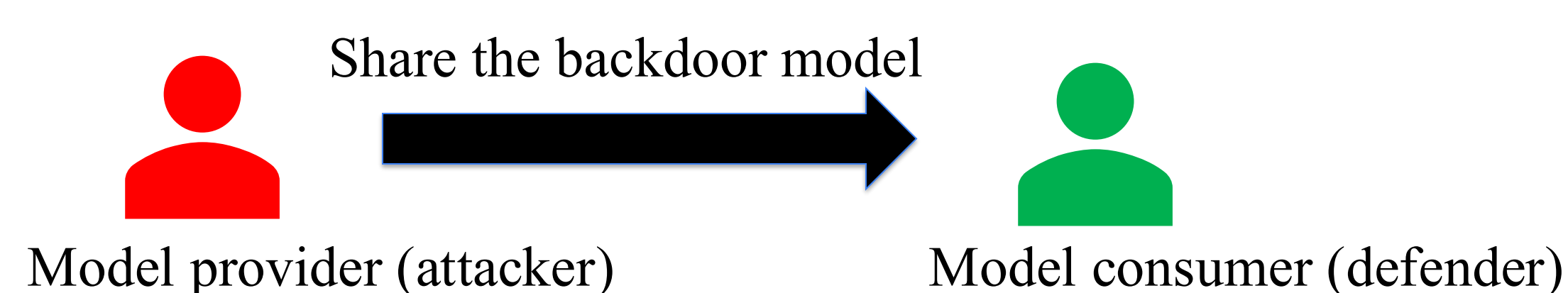


[2]



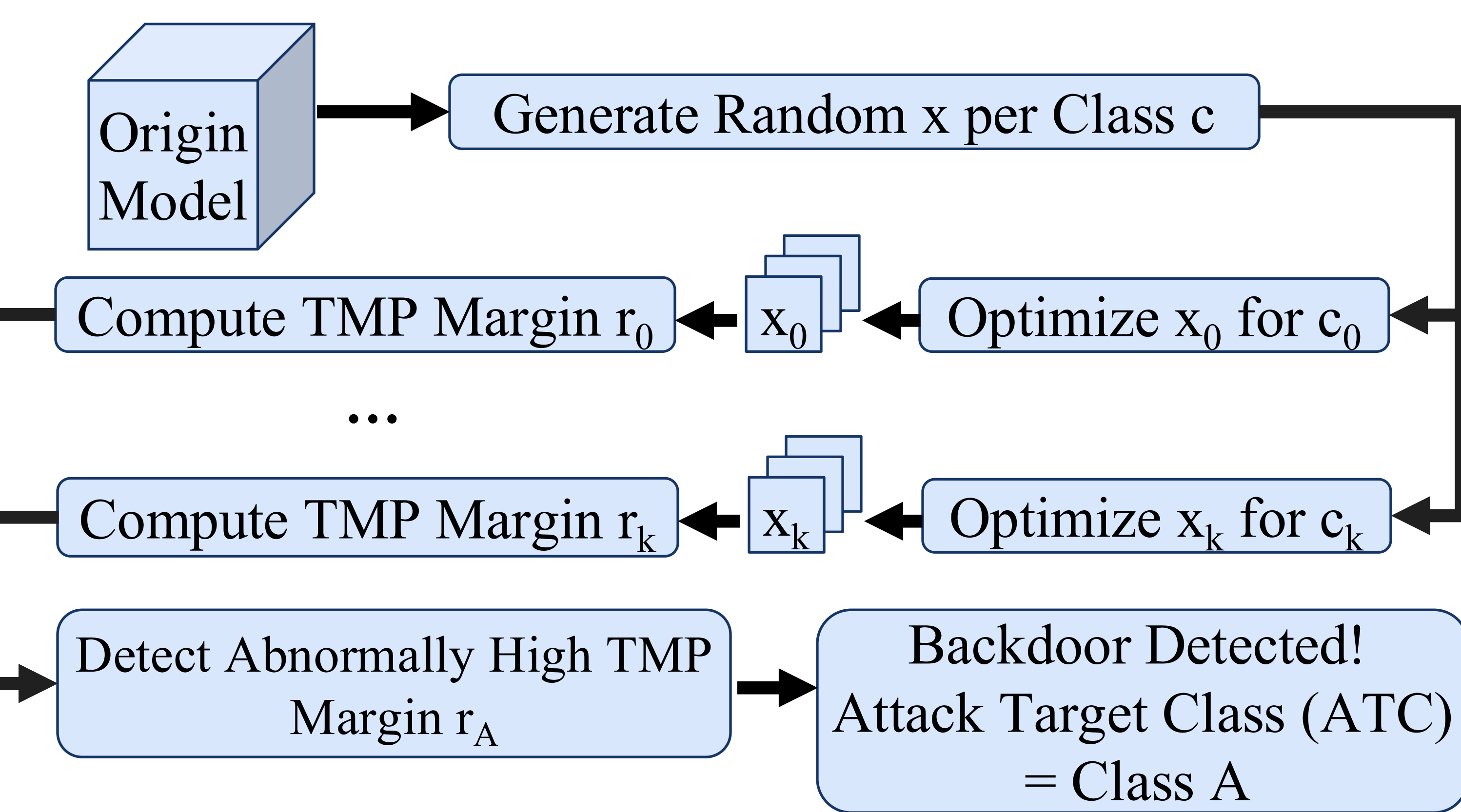
## Threat Model

- Attacker: dirty-label, all-to-one, data-poisoning backdoor attack.
- Defender: no training data, no clean labels, white-box model only.



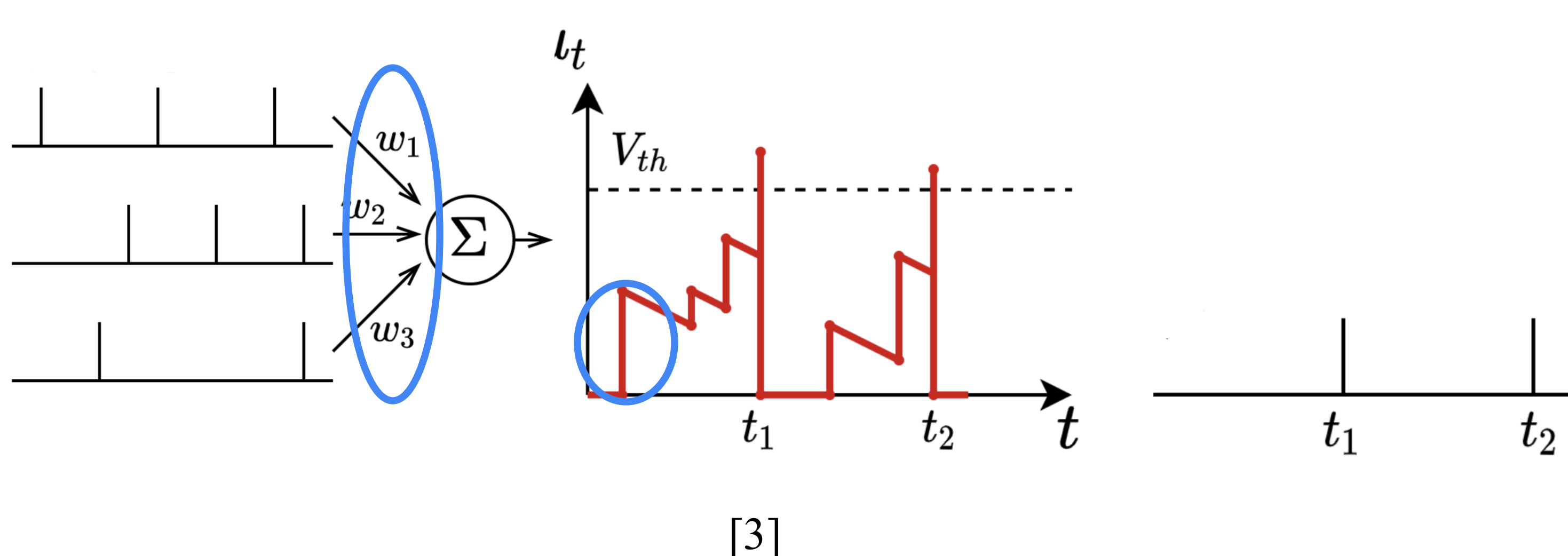
## TMPBD (Detection Block)

- Backdoor Attack Target Class (ATC) causes overactivation observed as abnormally high membrane potential.

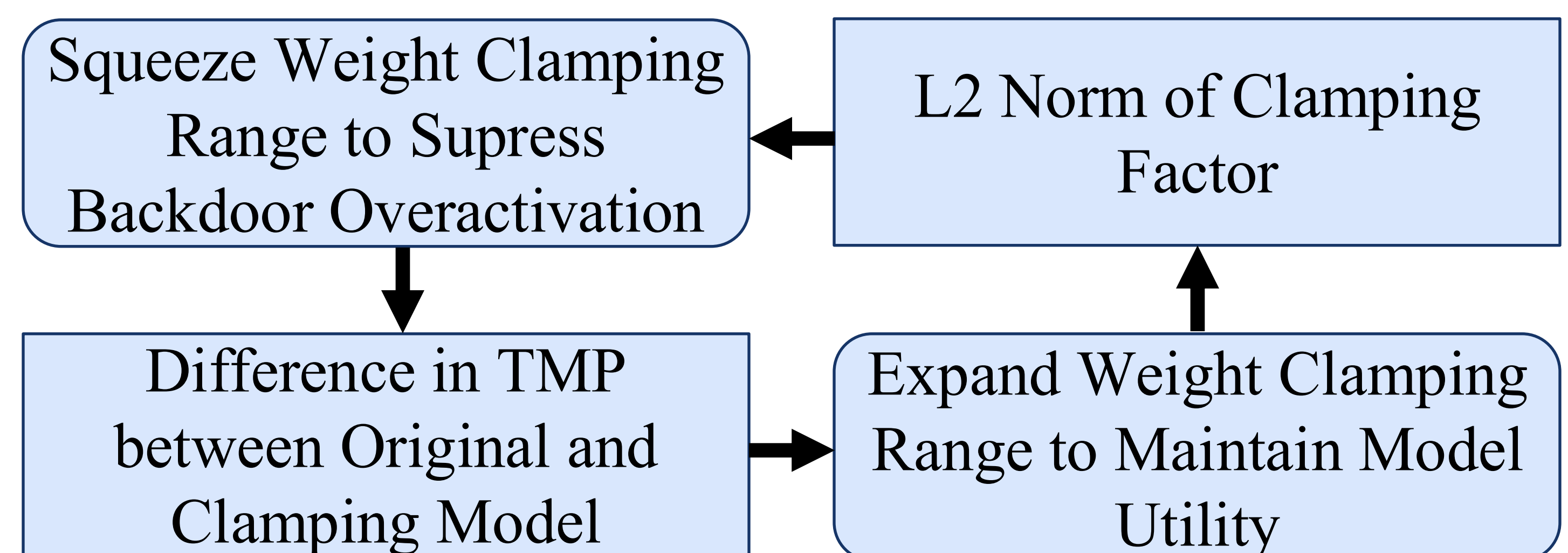


## NDSBM (Mitigation Block)

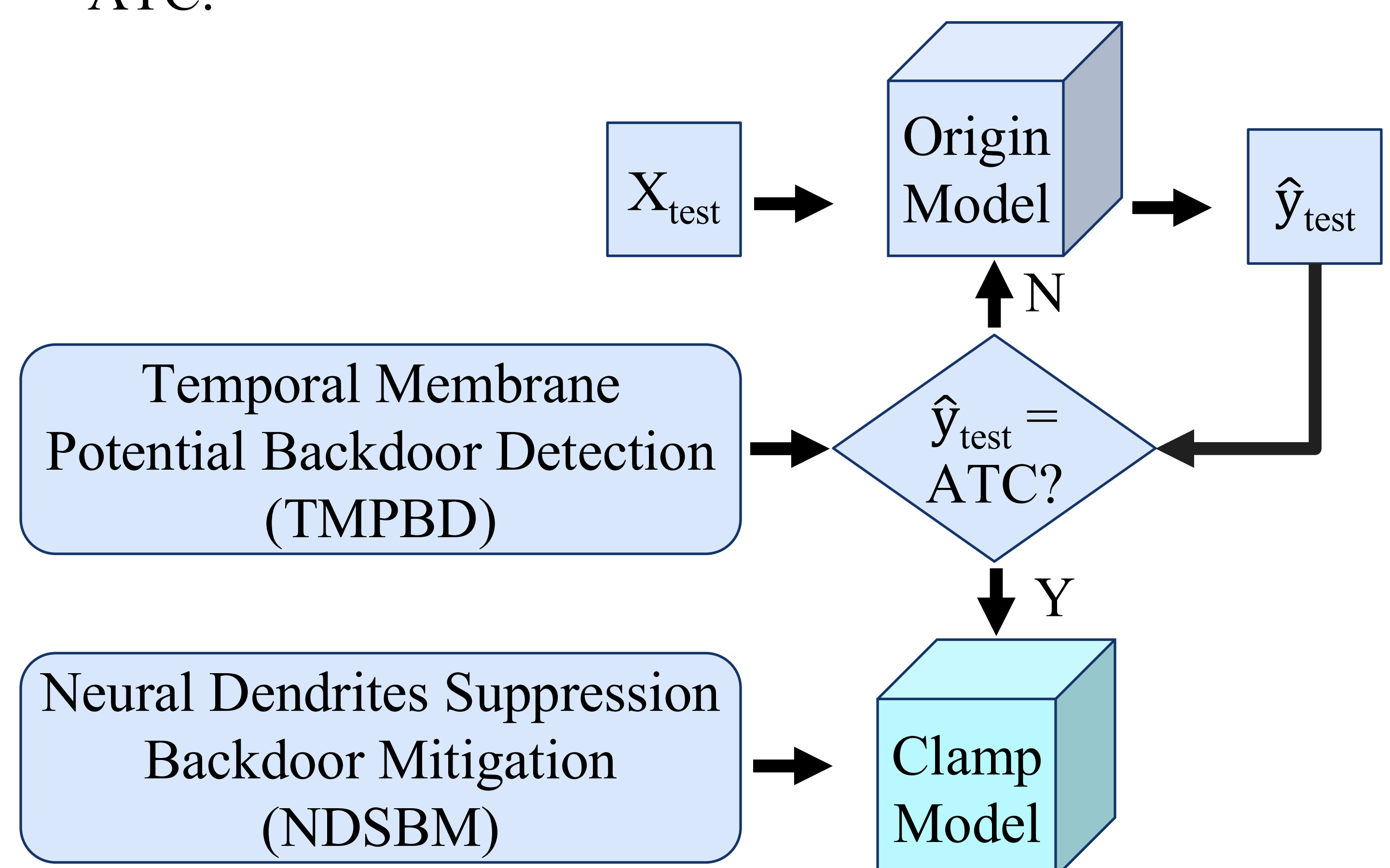
- Squeeze model weight into a clamping range to mitigate backdoor overactivation.



- Choosing weight clamping weight balancing the defense effectiveness and model utility.



- Only apply clamp model with inference sample predicted ATC.



## Results

- 100% backdoor & ATC detection accuracy.
- ASR: 100%  $\rightarrow$  2.81%.

	Clean		Backdoored	
	CA $\uparrow$	ASR $\downarrow$	CA $\uparrow$	ASR $\downarrow$
Original	97.65 $\pm$ 1.03	0.31 $\pm$ 0.99	84.71 $\pm$ 12.48	100.00 $\pm$ 0.00
Supervised				
Fine-Tuning	64.56 $\pm$ 6.63	3.00 $\pm$ 4.70	88.53 $\pm$ 5.54	3.28 $\pm$ 4.51
MMBM	73.09 $\pm$ 8.38	2.90 $\pm$ 2.99	71.76 $\pm$ 16.08	1.40 $\pm$ 2.49
Unsupervised				
Self-Tuning	7.20 $\pm$ 3.13	11.44 $\pm$ 19.89	6.47 $\pm$ 1.73	25.56 $\pm$ 20.18
Max Cla.	83.09 $\pm$ 3.81	2.28 $\pm$ 4.30	88.83 $\pm$ 4.00	19.38 $\pm$ 13.39
Abs. Cla.	84.26 $\pm$ 4.75	1.34 $\pm$ 2.09	89.12 $\pm$ 2.52	20.81 $\pm$ 14.50
NDSBM	72.50 $\pm$ 6.43	3.69 $\pm$ 10.27	89.86 $\pm$ 3.21	8.44 $\pm$ 9.91
<b>Ours</b>	<b>97.06<math>\pm</math>1.55</b>	<b>0.31<math>\pm</math>0.99</b>	<b>92.06<math>\pm</math>4.29</b>	<b>2.81<math>\pm</math>3.95</b>

## Reference

- [1] A. Viale, A. Marchisio, M. Martina, G. Masera, and M. Shafique, "Carsnn: An efficient spiking neural network for event-based autonomous cars on the loihi neuromorphic research processor," in *2021 IJCNN*, IEEE, 2021, pp. 1–10.
- [2] T. Gu, K. Liu, B. Dolan-Gavitt, and S. Garg, "BadNets: Evaluating Backdooring Attacks on Deep Neural Networks," *IEEE Access*, vol. 7, pp. 47230–47244, 2019, doi: [10.1109/ACCESS.2019.2909068](https://doi.org/10.1109/ACCESS.2019.2909068).
- [3] H. Kamata, Y. Mukuta, and T. Harada, "Fully spiking variational autoencoder," in *AAAI 2022*, pp. 7059–7067.



Full Paper



Source Code