PRACTICAL MECHANISM VIA SIMPLE INPUT CONTROL FOR FAULT-TOLERANT SPIKING NEURAL NETWORKS

Anonymous authorsPaper under double-blind review

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

025

026027028

029

031

032

033

034

037

040

041

042

043

044

046

047

048

049

050 051

052

ABSTRACT

Spiking Neural Networks (SNNs) attract researchers due to their energy-efficient operations in neuromorphic devices. Despite their energy efficiency, SNNs are vulnerable to hardware faults, which impair the functionality of learnable parameters (e.g., Stuck-At-Faults (SAFs) in synaptic weights). This impairment reduces the capacity to absorb information. When input data contains information exceeding the capacity, SNNs may not absorb information correctly, referred to as the **bottleneck problem**. Existing approaches have relied on complex algorithms or direct modification to most synaptic weights in SNNs, limiting their practicality in neuromorphic devices. This paper proposes a simple yet effective input control mechanism to address the problem, grounded in a thorough motivation study. Our mechanism divides the input samples into small fragments, following the best fragmentation strategy, derived by analyzing the characteristics of the input samples and diagnosing the current influence of faults. Experimental results demonstrate that our mechanism significantly enhances fault tolerance over existing methods, achieving these gains without complex algorithms or direct weight modification in various SNN models. Additionally, our mechanism improves the fault tolerance of SNN models implemented in a Field-Programmable Gate Array (FPGA) device.

1 Introduction

Spiking Neural Networks (SNNs) attract researchers to develop neuromorphic devices that are necessary to implement Artificial Intelligence (AI) in low-end devices (Garaffa et al., 2021; Jeong et al., 2025). SNNs are third-generation neural networks that use spikes to process data. They are well-suited for neuromorphic devices with limited power sources and biological operations because SNNs consume less energy than other neural networks and have high bio-plausibility (Schuman et al., 2022; Pfeiffer & Pfeil, 2018). Although SNNs are essential for neuromorphic devices, they remain highly vulnerable to permanent hardware faults, which frequently occur in the electrical components of neuromorphic devices and significantly impair SNNs' learning performance (Spyrou et al., 2021; Lee & Lim, 2023). Their fault vulnerability stems from the instability of neuromorphic devices' hardware components and the fault sensitivity of SNNs' neuron models (Garaffa et al., 2021).

Previous approaches to improve SNNs' tolerance against hardware faults rely on complex algorithms to regulate abnormal neuronal activities or demand hardware reconfigurability to manage electric components directly in neuromorphic devices (Vu et al., 2019; Putra et al., 2022). Although these approaches have improved SNNs' fault tolerance, they exhibit the following problems, which reduce their practicality in implementation.

- 1. Algorithmic complexity in conventional mechanisms: The previous approaches require complex algorithms that are impractical for hardware implementation (Vu et al., 2019; Yang et al., 2022; Han et al., 2023). These approaches cannot work properly in neuromorphic devices demanding low-power operations and prevent neuromorphic devices from operating stably (Basu et al., 2018). This is because the complex algorithms frequently malfunction due to unexpected events such as wrong input values and hardware faults (Liu et al., 2017).
- 2. Difficulties for direct modifications to synapses: Modifying synaptic weights directly is essential for the previous approaches to improve SNNs' fault tolerance, such as synapse pruning and weight bounding. These approaches forcibly adjust whole synaptic weights and configurations in neuromorphic devices (Putra et al., 2022; Chen & Chakrabarty, 2021). However, external methods enabling

direct modification demand additional costs to design reconfigurable hardware (Garaffa et al., 2021; Takano & Amano, 2022; Putra et al., 2023).

Pragmatic approaches to meeting these two problems are necessary to improve SNN's fault tolerance in the real world. To address these limitations, we identify a critical issue that significantly degrades SNN performance in faulty neuromorphic devices through a detailed motivation study. We name this issue **the bottleneck problem**, severely damaging SNNs' usable learning capacity. Here, we explain how the bottleneck problem occurs. When faults appear in SNNs' synapses, the weights of the faulty synapses become fixed during training. This means that the capacity is reduced because the faulty synapses, which do not change their weight during training, cannot be used for memorizing data. Furthermore, the pre-activation value (linear combination) of spiking neurons lies in an abnormal point of the surrogate gradient function, causing a serious gradient vanishing problem, which results in the capacity gradation. With the low usable capacity of faulty SNNs, they cannot memorize input data properly when the data contains information exceeding the learning capacity.

Motivated by flow control methods used in computer networks (Kurose & Ross, 2012), we first propose a practical mechanism based on simple input control to solve the bottleneck problem based on data fragmentation. Our mechanism enhances the fault tolerance of SNNs by dividing input data samples into small fragments. Our control scheme addresses the constraint of SNNs' usable learning capacity in faulty neuromorphic devices by exploiting an effective fragmentation strategy for fault mitigation based on analysis of input images' characteristics and the influence of faults. The novelty of our mechanism is as follows. Unlike the previous approaches, our mechanism does not require complex algorithms to control neuronal activities and hardware reconfigurability for modifying synapses directly. This novelty is derived from the following features of our mechanism.

- 1. Data sample division into small fragments: Our algorithm fragments input data samples to decrease the input samples' size and shrink information in the input samples, leading the faulty SNN models to memorize the information despite their degraded usable learning capacity due to faults.
- 2. Fragmentation method to minimize the adverse effect from faults: To ensure that our mechanism provides the fragments that models can handle, we develop a fragmentation strategy that adapts fragment geometry (i.e., the cut angle) to mitigate fault-induced damage during the forward pass.

We develop our mechanism through a thorough motivation study with faulty SNN models. With our mechanism, various SNN models achieve significantly higher classification accuracy than models using previous approaches for fault mitigation under fault-injected conditions, while consuming less energy due to our simple approach. We additionally conduct experiments with real hardware SNNs built in a Field-Programmable Gate Array (FPGA) device. Our work has the following contributions.

- We develop a practical mechanism to enhance SNNs' fault tolerance without complex algorithms and direct synapse modifications, which easily malfunction and demand high hardware reconfigurability.
- We present a concrete theoretical basis for our mechanism by mathematically and experimentally investigating how synaptic faults degrade the usable learning capacity of SNN models in our motivation study. Due to our detailed motivation study, we propose a theoretically sound fragmentation mechanism.
- We provide a rich set of evaluation results in various scenarios, including Hardware SNNs based on an FPGA device. The evaluation results demonstrate that our mechanism exhibits better fault mitigation ability than previous approaches by using simple methods to enhance the fault tolerance of SNNs.

2 BACKGROUNDS

2.1 Spiking neural networks

SNNs are third-generation neural networks motivated by the learning mechanisms of the human brain (Yao et al., 2023). In SNNs, spiking neurons fire and emit output spikes at every time interval, corresponding to each time step. Here, the time step is a unit of time for spike occurrence. The spiking neurons generate spikes only when the membrane potential of the neurons reaches a threshold.

Researchers use various neuron models to build SNN models. Among them, the Integrate-and-Fire (IF) model and Leaky Integrate-and-Fire (LIF) model, which have leakage in membrane potential unlike IF, are widely used (Moitra et al., 2023). Synaptic weights determine how significantly input spikes from pre-synaptic neurons of the previous layer affect the post-synaptic neurons (Venkatesha et al., 2021). SNNs update the synaptic weights with the two approaches: supervised and unsupervised learning rules. Supervised learning rules calculate gradients of spiking neurons to update weights using surrogate gradient functions. Unsupervised learning rules utilize the time difference between pre-synaptic and post-synaptic activity to update weights.

SNNs are necessary to implement neuromorphic devices. They exhibit less energy consumption than conventional neural networks for the following reasons. First, spiking neurons fire and update synaptic weights only when a specific event occurs (Lee & Lim, 2024). The SNN's infrequent spike generation is associated with sporadic data processing, resulting in low power consumption. Second, SNNs replace complex Multiply-ACcumulate (MAC) operations with simple ACcumulate (AC) operations, eliminating weight multiplication to input data while accumulating input information.

2.2 SYNAPTIC FAULTS

Synaptic faults are persistent or transient defects in the synaptic weights of a connection. They distort pre-activations and inject biased or structured noise, which disrupts the learning process of SNN models. A representative case of them is the Stuck-At Faults (SAFs), where a weight is fixed to the highest (SA1) and lowest (SA0) synaptic weight, ignoring updates and introducing systematic bias (Vatajelu et al., 2019). Note that SA0 and SA1 do not mean the weights are stuck at the exact values 0 and 1. Another common case is Random Weight Faults (RWFs), transiently removing intended synaptic connections or creating unintended connections due to thermal noise (Vatajelu et al., 2019). Additionally, Connectivity Error Faults (CEFs) permanently change the connections between synapses, ruining synaptic connections. These faults shift the pre-activation away from useful operating regions, shrink the effective gradient signal, and reduce usable learning capacity.

3 STATE OF THE ARTS

3.1 Analysis about faults in neuromorphic devices

Researchers have deeply investigated how faults affect neuromorphic devices. They inject faults into synapses and neurons of SNNs in neuromorphic devices, and how these faults ruin the classification performance of SNNs (Vatajelu et al., 2019). They build a memristive neuromorphic simulator and analyze how faults disturb data classification (Lee & Lim, 2023). In this study, researchers prove that the faults occurring in synapses correlated with important features of data samples influence the devices more severely. Researchers also apply various fault types in neuromorphic devices and study how spiking neurons act in detail (Ali El Sayed, 2021; Garaffa et al., 2021). However, these analysis studies overlook the excessive updates caused by hardware faults in neuromorphic devices.

3.2 MECHANISMS TO IMPROVE FAULT TOLERANCE OF SNNS

The conventional methods to improve the fault tolerance of SNNs are based on dropping faulty elements in SNNs and applying additional architectures for fault mitigation. Researchers utilize the error correction ability of binary codes in output decoding to enhance the fault tolerance of neural networks (Liu et al., 2019; Yu et al., 2023). They also induce spikes to avoid faults in SNNs to reduce the bad effects caused by faults (Vu et al., 2019; Yang et al., 2022). They build a fault map to select spiking neurons severely affected by faults and shrink these neurons' influences (Putra et al., 2022; Wicaksana Putra et al., 2021; Yang et al., 2022). Additionally, they mask faulty elements by setting affected pre-trained weights to zero, then retrain with per-layer threshold (Siddique & Hoque, 2023). Researchers also employ self-recovering mechanisms from astrocytes that recover sick neurons in the human brain to neuromorphic devices. With the self-recovering ability of astrocytes, neuromorphic devices' fault tolerance improves (Han et al., 2023; Varshika et al., 2023). Enhancing the astrocyte-based approaches, they augment SNNs with an astrocyte-inspired leaky integrator, stabilizing spiking dynamics and markedly improving fault tolerance (Yunusoglu et al., 2025). Despite their enhancement of fault tolerance, they require complex architectures based on complicated algorithms and neglect the reconfigurability of electric components in hardware neuromorphic devices.

4 MOTIVATION STUDY

4.1 OVERVIEW

We discovered that the faults cause the bottleneck problem with the following procedures.

- 1. Synaptic faults increase the pre-activation magnitude |z|. The pre-activation is a linear combination of inputs: in fully-connected layers: z=Wx+b (component-wise $z_j=\sum_i w_{ji}x_i+b_j$); in convolution layers: $z_o(u,v)=\sum_c\sum_{\Delta u,\Delta v}w_{o,c,\Delta u,\Delta v}\,x_c(u+\Delta u,v+\Delta v)+b_o$. When faults perturb the weights to $W+\Delta W$ (or $w+\Delta w$), the pre-activation changes significantly.
- 2. As |z| grows and moves away from the (spiking) threshold of spiking neurons, the surrogate gradient near the threshold collapses toward zero, so the learning signal cannot propagate backward.
- 3. The faulty weights are fixed to abnormal values, and the non-faulty weights barely change due to near-zero gradients. This problem significantly reduces SNNs' usable learning capacity, creating a bottleneck that prevents the model from fitting the data it should learn.

We mathematically explain how faults fail the learning process of SNN models in Appendix C.

4.2 PRE-ACTIVATION MAGNITUDE INCREASE BY FAULTS

To demonstrate that the synaptic faults cause the pre-activation magnitude to increase, we inject the SAFs (SA0:SA1=1.75:9.04) (Chen et al., 2017) and RWFs (representative permanent and transient faults) into 50% of synapses in a spiking Multi-Layered Perceptron (MLP) with 4 layers, VGG-7, and ResNet-18 with MNIST, CIFAR-10, and CIFAR-100, observing the change in pre-activation magnitude of all neurons in the SNN models according to faults after training with an Adam optimizer and Root-Mean-Square Error (RMSE) as a loss function. We obtain the experimental results by repeating experiments 10 times and present the pre-activation magnitude as a 95% confidence interval.

Table 1: Pre-activation magnitude of MLP, VGG-7, and ResNet-18 SNN models under fault injection.

	MLP(MNIST)	VGG-7(CIFAR-10)	RESNET-18(CIFAR-100)
NOMINAL SAFS RWFS	$\begin{array}{c} 2.7 \pm 0.13 \\ 393.42 \pm 10.89 \\ 8.08 \pm 1.75 \end{array}$	8.75 ± 2.08 273.66 ± 13.69 205.73 ± 10.27	$\begin{array}{c} 14.54 \pm 3.56 \\ 140.03 \pm 15.77 \\ 103.41 \pm 11.8 \end{array}$

Table 1 compares the summation of the pre-activation magnitude of all spiking neurons in the MLP (LIF), VGG-7 (LIF), and ResNet-18 (IF) SNN models with and without SAF and RWF injection. Our experimental results show that SAFs and RWFs significantly increase the pre-activation magnitude around all neurons with fault-injected synapses.

4.3 GRADIENT COLLAPSE BY ABNORMAL PRE-ACTIVATION MAGNITUDE

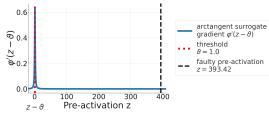


Figure 1: The value (average on the layers of the model) of surrogate gradient $\phi'(z-\vartheta)$ (arctangent) when we inject SAFs into the synapses of the MLP model.

We call the region where the surrogate derivative is non-negligible the surrogate gradient corridor and denote its half-width by δ (or threshold-aligned bound z^\star); outside this corridor, $\phi'(z-\vartheta)\approx 0$. As the pre-activation magnitude increases, the pre-activation value moves away from the corridor. This alignment error, due to an abnormal pre-activation magnitude, makes the surrogate gradient values nearly zero.

Figure 1 depicts the surrogate gradient function (arctangent) of an LIF neuron model and the position of pre-activation of the SAF-injected neurons in the MLP model. The surrogate gradient $\phi'(z-\vartheta)$ of

the SAF-injected neurons is near zero, and the gradient of these neurons vanishes. This is because the following equation calculates the gradient: $\delta^{(l)} = (\nabla_{a^{(l)}} L) \odot \phi'(z^{(l)} - \vartheta)$ (L is a loss function and l is the index of a layer).

219 220 221

4.4 LEARNING ABILITY DEGRADATION BY GRADIENT COLLAPSE

223 224 225

222

We demonstrate that gradient vanishing due to faults causes degradation in the SNNs' learning ability. Table 2: The gradients' L1 norm in 95% confidence interval upon all neurons and classification accuracy of MLP, VGG-7, and ResNet-18 SNN models under fault injection with 50% fault ratio.

226 227

228 229 230

	MLP(MNIST)	VGG-7(CIFAR-10)	RESNET-18(CIFAR-100)
NOMINAL (L1 NORM)	11.67 ± 1.71	25.68 ± 2.89	7.76 ± 1.92
NOMINAL (ACCURACY)	$97.57 \pm 0.38\%$	$56.63 \pm 0.91\%$	$25.59 \pm 1.24\%$
SAFs (L1 NORM)	0.01 ± 0.00026	1.83 ± 0.16	1.36 ± 0.25
SAFs (ACCURACY)	$11.35 \pm 0.01\%$	$9.99 \pm 0.08\%$	$5.04 \pm 1.07\%$
RWFs (L1 NORM)	65.86 ± 10.57	2.18 ± 0.45	0.49 ± 0.07
RWFs (ACCURACY)	$79.8\pm8.74\%$	$22.28 \pm 5.34\%$	$1.5 \pm 0.51\%$

231 232 233

234

235

236

237

238

Table 2 shows the gradients' L1 norm upon all neurons in the SNN models during training and the classification accuracy of the models with the various datasets after testing. The classification accuracy is proportional to the L1 norm of the gradient under fault injection with 50% fault ratio. This point demonstrates that vanishing surrogate gradients stall weight updates, preventing the network from making new decision boundaries and thus directly reducing its learning capacity (occurrence of the bottleneck problem). Interestingly, the MLP model does not accurately classify MNIST data samples despite a large gradient L1 norm when the model is under RWF injection. This occurs because RWFs induce random changes in synaptic weights, which transiently perturb the loss and cause a temporary increase in the gradient (Foret et al., 2021). Since the MLP model has limited learning ability, it cannot effectively compensate for such perturbations.

4.5 PRE-ACTIVATION SENSITIVITY OF LAYERS IN MLP

243 244 245

Our motivation study outlines an interesting finding: in MLP, the gradients of neurons in the previous layers are more sensitive to abnormal changes in pre-activation by faults than the later layers.

250 251

253

254

255

256

Table 3: The gradients' L1 norm in 95% confidence interval upon all neurons in each layer of the MLP model under fault injection.

	Layer 1	Layer 2	Layer 3	Layer 4
MNIST (SAFs)	$(3.18 \pm 1.13) \times 10^{-6}$	$(3.31 \pm 1.64) \times 10^{-3}$	$(1.15 \pm 0.93) \times 10^{-3}$	$(2.67 \pm 0.65) \times 10^{-3}$
MNIST (RWFs)	$(5.61 \pm 1.46) \times 10^{-5}$	$(5.19 \pm 1.52) \times 10^{-5}$	$(1.37 \pm 0.58) \times 10^{-4}$	$(1.39 \pm 0.29) \times 10^{-3}$
FMNIST (SAFs)	$(3.45 \pm 0.91) \times 10^{-15}$	$(1.01 \pm 0.37) \times 10^{-12}$	$(1.6 \pm 0.69) \times 10^{-7}$	$(1.02 \pm 0.36) \times 10^{-2}$
FMNIST (RWFs)	$(4.91 \pm 1.55) \times 10^{-5}$	$(5.11 \pm 1.71) \times 10^{-5}$	$(1.46 \pm 0.62) \times 10^{-4}$	$(1.76 \pm 0.82) \times 10^{-3}$

257 258 259

Table 3 shows the gradients' L1 norm of all neurons in each layer of the MLP model. The error signal that reaches layer ℓ is obtained by repeatedly applying the Jacobians of all higher layers in MLP. Under faults, pre-activations drift away from the operating threshold, so the surrogate derivatives $\phi'(z-\vartheta)$ on affected layers become very small. During gradient calculation, gradients are multiplicatively contracted by a chain of small factors. Because earlier layers (smaller ℓ) accumulate more of these factors, they suffer disproportionately severe gradient vanishing, explaining the front-loaded degradation we observe under faults.

260 261

4.6 SIMILARITY TO FLOW CONTROL IN COMPUTER NETWORKS

267

268

269

Data flow control mitigates the congestion problem in computer networks (Kurose & Ross, 2012; Wigren & Karaki, 2018). The information of input samples in SNNs is related to the data in the packets of computer networks, and the surrogate gradient corridor is related to the data capacity that a receiver can handle in computer networks. The bottleneck problem in neural networks is also similar to that of the receiver, which prevents the receiver from processing the large data packet simultaneously. The object of flow control, adjusting the size of data in a packet to satisfy the receiver's data capacity, is similar to the object of enhancing fault tolerance: changing the size of information in the input samples to keep the pre-activation in the surrogate gradient corridor.

5 PROPOSED MECHANISM

Our motivation study demonstrates that synaptic faults can inflate pre-activations beyond the surrogate gradient corridor, causing the gradient collapse. To prevent the gradients from vanishing, we design an adaptive input fragmentation mechanism to avoid drift in pre-activation magnitude by shrinking the probability of cases where input data samples enter faulty synapses, which have abnormal weights that cause the pre-activation magnitude to increase significantly. We mathematically prove why our mechanism is nearly-optimal in Appendix D.

5.1 Sensitivity score definition and calculation

Key point 1: The sensitivity score represents which pixel changes the pre-activation the most significantly under faults.

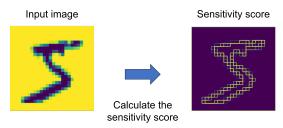


Figure 2: Sensitivity score calculation with an MNIST image.

Figure 2 depicts the sensitivity score of an MNIST image. To measure how input data samples are affected by synaptic faults and explore the best fragment shape that prevents a significant increase in pre-activation magnitude, we define the sensitivity score of input samples, consisting of an image sample's complexity and the <u>influence</u> measurement from faults.

$$I(x) = PN\Big((LoG(x) + Sobel(x) + Var(x)) \odot (1 + \lambda_s SP(x) + \lambda_w WP(x)) \Big).$$
 (1)

where $\mathrm{LoG}(x)$, $\mathrm{Sobel}(x)$, and $\mathrm{Var}(x)$ are information about edges, blobs, and texture contrast of images. We adopt them to calculate the complexity of an image sample and design our algorithm to minimize each fragment's complexity (Lowe, 2004). This is because an image with high complexity significantly changes the pre-activation since its norm is large. $\mathrm{SP}(x)$ is a saliency map $\propto \|\partial \mathcal{L}/\partial x\|_{2,\mathrm{channel}}$ and $\mathrm{WP}(x)$ is the absolute first layer's weight projection to input resolution. They convert the pixel-derived map into a fault-influence field on the weights of the first layer. We apply it to our mechanism to minimize the increase in pre-activation by reducing the probability that pixel values enter many faulty synapses at once. PN is a percentile normalization, which normalizes the values of the fault influence map in the range of 0 to 1. $I(x) \in [0,1]^{H \times W}$. For batch stability and fast operation, we measure the sensitivity score of the averaged image sample of a batch.

5.2 GINI COEFFICIENT CALCULATION WITH A 1D PROFILE

Key point 2: The Gini coefficient indicates the equality of the sensitivity score. It should be minimized to increase the equality of the sensitivity score upon the fragments and prevent the pre-activation from falling outside of the surrogate gradient corridor.

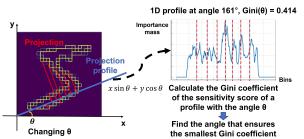


Figure 3: 1D projection and Gini coefficient calculation changing the projection angle.

As depicted in figure 3, changing the angle θ of the 1D profile by 1° and projecting the pixels on the profile with the equation $s = xsin(\theta) + ycos(\theta)$, we calculate the Gini coefficient with equation 2

$$Gini(p_{\theta}) = \frac{1}{2L\bar{p}_{\theta}} \sum_{i=0}^{L-1} \sum_{j=0}^{L-1} |p_{\theta}[i] - p_{\theta}[j]|.$$
 (2)

where, $p_{\theta} \in \mathbb{R}_{\geq 0}^{L}$ is the 1D importance profile obtained by projecting along angle θ (range of $[0^{\circ}, 360^{\circ}]$). $p_{\theta}[\bar{i}]$ is the value of bin i ($i = 0, \ldots, L-1$). L is the number of bins and $\bar{p}_{\theta} = \frac{1}{L} \sum_{k=0}^{L-1} p_{\theta}[k]$ is the mean. $|\cdot|$ denotes absolute value. The factor $(2L\bar{p}_{\theta})^{-1}$ normalizes the sum of pairwise absolute differences, so the index is 0 for a uniform profile and grows with inequality.

We explore the angle θ which makes the 1D profile have the minimum Gini coefficient. This is because the Gini coefficient is a strictly Schur-convex function, ensuring the Gini coefficient strictly increases under majorization: if $x \succ y$ (i.e., x is more unequal), then $\mathrm{Gini}(x) > \mathrm{Gini}(y)$, with equality only for permutations (Sandor, 2007). By minimizing the Gini coefficient, the fragment's equality of sensitivity score is maximized. We design our mechanism to maximize equality and prevent the pre-activation from leaving the corridor. This is because as a fragment's sensitivity score becomes more equal by minimizing Gini coefficient, the maximum of each fragment's energy $\|x_t\|_2$ decreases, so the upper bounds $\|w^\top x_t\| \le \|w\|_2 \|x_t\|_2$ and $\|u_t\| = \|w^\top x_t + m\|$, $(m = b - \vartheta)$ are pushed below δ , keeping the pre-activation fixed inside the corridor (b is a bias in SNNs' layers.).

5.3 Fragment generation based on equal sensitivity score and in/out for SNNs

Key point 3: The fragmentation line is set by the 1D profile cutting to make the 1D bins have an equal cumulative sum of the sensitivity score.

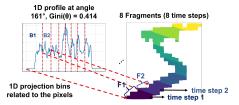


Figure 4: Dividing an image sample based on the cumulative sum of the sensitivity score.

After finding the angle that ensures the minimum Gini coefficient, we calculate the cumulative sum of the sensitivity score. As depicted in figure 4, the division line for fragmentation (dashed line in the figure) indicates the position to cut the 1D profile into bins B_1, \ldots, B_T , making the bins have equal sensitivity score. We generate fragments $F_1, \ldots, F_T \in \{0,1\}^{H \times W}$ with the pixels correlated to the points in bins B_1, \ldots, B_T and feed the SNN models the fragments F_1, \ldots, F_T over time steps t=1..T. We zero-pad the generated fragment to align the input dimension of fragments with the input dimension of the original samples because the input dimension of SNNs is not changeable during training and testing. We apply the same division angle to all data samples in the same batch.

To ensure that pre-activation z_t of fragment F_t positions in the corridor, keeping the scale of active pixels in each fragment to a target Root-Mean-Square (RMS) is also important. Therefore, we adopt RMS normalization to the input fragments with the equation 3.

$$\tilde{x}_t = g_t x_t, \quad \|\tilde{x}_t\|_2 = \alpha \implies |z_t| \le \|\hat{w}\|_2 \alpha. \tag{3}$$

where α is the non-zero pixels in the fragment for the input of time step t. $\|\hat{w}\|_2 \alpha \leq z_{\star}$ places z_t inside the surrogate-derivative corridor.

We also adopt an entropy-based output decoding technique to aggregate the outputs (logits) from SNNs across all time steps accurately.

$$\bar{\ell} = \sum_{t=1}^{T} e_t \,\ell_t, \qquad e_t \propto \exp\{-\tau H(\operatorname{softmax}(\ell_t))\}.$$
(4)

where $\ell_t \in \mathbb{R}^C$ is an output vector from SNNs at time t (with C classes) and T is the number of steps. Letting $p_t = \operatorname{softmax}(\ell_t)$, the Shannon entropy is $H(p_t) = -\sum_{c=1}^C p_{t,c} \log p_{t,c}$ and $\tau \geq 0$ controls how strongly low-entropy (confident) steps are emphasized. The scale factor e_t to represent

the entropy is $e_t \propto \exp\{-\tau H(p_t)\}$ (where \propto denotes proportionality and the e_t is normalized so $\sum_{t=1}^T e_t = 1$), and the entropy-weighted aggregate logit is $\bar{\ell} = \sum_{t=1}^T w_t \, \ell_t$ for final decoding.

6 EXPERIMENTS

6.1 Experimental settings

We conduct various experiments with MLP (LIF neurons), VGG-7/11/15 (LIF neurons), and ResNet-18/34 (IF neurons) SNN models based on SpikingJelly, widely used for SNN implementation, by classifying samples in MNIST/FMNIST/UCI-HAR (MLP), CIFAR-10/100 (VGG and ResNet), and Tiny-ImageNet (ResNet) (Fang et al., 2023). We select these SNN models and datasets because current SNN technologies do not appropriately train large and deep models with complicated datasets (Fang et al., 2023; Schuman et al., 2022). We measure the classification accuracy of the SNN models using the proposed mechanism and benchmarks under SAFs, setting the ratio of SA1 and SA0 to SA0: SA1 = 1.75: 9.04 (Chen et al., 2017) and the weight boundary to [-1, 1] (Le Gallo et al., 2023; Lammie et al., 2022). Additionally, we inject RWFs and CEFs into the synapses of these models. We use ECOC (Liu et al., 2019), SoftSNN (Putra et al., 2022), Routing (Vu et al., 2019), Astrocyte (Han et al., 2023), FalVolt (Siddique & Hoque, 2023), and LIFA (Yunusoglu et al., 2025) for our benchmarks¹. The SNN models without any fault mitigation mechanism are the baseline. We use RMSE as a loss function and Adam as the optimizer for SNN models (Fang et al., 2023). The batch size is 100, and the learning rate is 0.001 for MLP/VGG and 0.01 for ResNet. We set the number of time steps (fragments) to 2, 4, and 8. We use 50 epochs for training. We set λ_s , λ_w to 0.1, and τ to 2.0 by tuning these parameters through experimental repetition with the proposed mechanism. We repeat all experiments 10 times with different random seeds and present the experimental results in a 95% confidence interval. We inject faults into the synapses sporadically in a uniform distribution, resulting in the uniform position of synaptic faults. Note that the additional results from the additional datasets (UCI-HAR and Tiny-ImageNet), different time steps (4 and 8 steps), other fault types (RWFs and CEFs), the ablation study with the combination of our mechanism, various hyperparameter (λ_s and λ_w) settings, and evaluations with an actual FPGA device, are presented in Appendix A. Additionally, we show the results with DNN models in Appendix A.

6.2 CLASSIFICATION ACCURACY COMPARISON

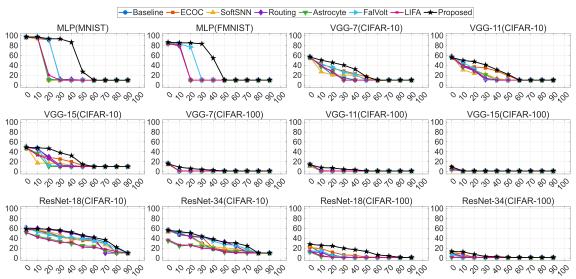


Figure 5: Average classification accuracy of various SNN models with the baseline, benchmarks, and proposed mechanism under SAFs using 2 time steps. The x-axis is the fault ratio (%) and the y-axis is the accuracy (%).

Figure 5 illustrates the classification accuracy of SNN models compared to the baseline, benchmarks, and proposed mechanism under SAFs with 2 time steps. Under SAFs, the SNN models with our mechanism exhibit the best classification accuracy across all datasets and models in most cases.

¹We briefly explain how these benchmarks enhance the fault tolerance in Section 3.

6.2.1 MLP MODELS

In the MLP model, the classification accuracy drops dramatically as the fault ratio increases. This is because the faulty weights are directly multiplied by the input values, and the pre-activation magnitude increases significantly, allowing it to easily escape from the surrogate gradient corridor. Our mechanism definitely outperforms the baseline and benchmarks, since it utilizes the input saliency and weight projection map of the first layer in Gini-based equal fragmentation. The MLP model is vulnerable to faults in the first layer, as mentioned in Subsection 4.1. By adopting the saliency and weight projection map (fault influence map) of the first layer, we suppress its pre-activation, which decides the surrogate gradient, from increasing significantly and escaping from the corridor. Thus, the pre-activation does not lie far from the corridor, preserving the power of the first layer's gradient.

6.2.2 VGG MODELS

While VGG models with the benchmarks using CIFAR-10 maintain proper classification performance only up to a fault ratio of 30–40%, the models with the proposed mechanism sustain correct classification even at fault ratios as high as 50%. This is because our mechanism targets to make the pre-activation lie at the point in the surrogate gradient corridor with Gini-based equal mass fragmentation, despite the large amount of faults. Contrarily, the benchmarks do not consider the relations between pre-activation, corridor, and surrogate gradient, failing to bound the pre-activation in the corridor. Therefore, the classification accuracy of the model using the benchmarks degrades sharply with fewer faults. We also observe that the classification accuracy declines as the model gets deeper. This occurs because surrogate gradients in deep SNNs cannot reliably approximate the hypothetical gradients of LIF neurons (Guo et al., 2024). When we use CIFAR-100, the models do not classify the data samples accurately due to their low learning ability. However, the models with our mechanism exhibit the highest classification accuracy under SAFs in most cases.

6.2.3 RESNET MODELS

Different from VGG models, ResNet models integrated with the benchmarks and proposed mechanism maintain the classification accuracy up to a fault ratio of 80-90% when we use CIFAR-10. This is because ResNet models have internal mechanisms to compensate for errors in gradient calculations, such as residual blocks. They also classify CIFAR-100 samples more accurately than VGG models under faults, since they have a more powerful learning ability than VGG models. ResNet-18 using CIFAR-100 maintains its classification ability up to the ratio of 60-70% only with our mechanism, and the ResNet-34 with CIFAR-100 maintains the classification ability up to the ratio of 30-40% only with our mechanism. These results demonstrate that our mechanism successfully enhances the fault tolerance with complicated datasets and models. We observe that the astrocyte-based approaches (Astrocyte and LIFA) do not improve the deep ResNet models' fault tolerance at all. This problem derives from the fact that they only mimic biological mechanisms of neuronal activity in brains, which enhances the fault tolerance of shallow and highly bio-plausible models such as Diehl & Cook 2015, using a bio-plausible unsupervised learning rule (Han et al., 2023; Yunusoglu et al., 2025). However, our mechanism successfully strengthens the models' fault tolerance in most cases since we tackle a fundamental problem of faulty SNN models regardless of the types of SNN models, and develop a solution to mitigate the problem.

7 CONCLUSION

This paper introduces a simple yet effective fault mitigation mechanism for SNNs that does not require complicated architectures or direct weight modifications based on input data control. Our mechanism improves fault tolerance more effectively than conventional approaches in various SNN models and datasets. Experimental results exhibited improvement in the fault tolerance of our mechanism over benchmarks in various network models and datasets, including real hardware environments. We emphasize that this improvement is primarily achieved through an effective input data control mechanism based on detailed observation of how synaptic faults ruin the learning capability of SNNs. Our mechanism allows SNNs to maintain reliable operation and high fault tolerance in a practical and hardware-compatible manner, enabling more sustainable and reliable edge AI computing.

ETHICS STATEMENT

This work does not involve human subjects, sensitive personal data, or potentially harmful applications. We trained and validated our models using publicly available datasets (e.g., MNIST, FMNIST, CIFAR-10, and CIFAR-100), without any private or identifiable information. We design our mechanism to enhance the robustness of neuromorphic systems against hardware faults. We declare no conflict of interest or external sponsorship that might have influenced the research outcomes.

REPRODUCIBILITY STATEMENT

We conducted all experiments on publicly available datasets with standard train/validation/test splits. To facilitate replication, we provide our full implementation of the proposed mechanisms in the anonymous supplementary material. We also present the settings of hyperparameters, model configurations, and hardware specifications to support reproducibility. We will publicly release the code and scripts on GitHub if our paper is accepted for the conference.

LARGE LANGUAGE MODEL USAGE STATEMENT

We used Large Language Models (LLMs) as writing and experiment assistants to improve the clarity of writing/editing mathematical equations (fixing typos/suggesting algebraic simplifications), and assist us in conducting experiments with the benchmarks. We did not use LLMs for idea generation, methodological design, analysis, or to originate any mathematical arguments or claims. We derived, verified, and finalized all derivations, results, and claims.

REFERENCES

- Sarah Ali El Sayed. Fault Tolerance in Hardware Spiking Neural Networks. Theses, Sorbonne Université, October 2021. URL https://theses.hal.science/tel-03681910.
- Arindam Basu, Jyotibdha Acharya, Tanay Karnik, Huichu Liu, Hai Li, Jae-Sun Seo, and Chang Song. Low-power, adaptive neuromorphic systems: Recent progress and future directions. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, 8(1):6–27, 2018. doi: 10.1109/JETCAS. 2018.2816339.
- Léon Bottou, Frank E. Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning, 2018. URL https://arxiv.org/abs/1606.04838.
- Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, Cambridge, 2004. doi: 10.1017/CBO9780511804441.
- BrainChipInc. Akd1000 akida system-on-chip product brief. Technical Report V2.3, BrainChip-Inc., 8 2025. URL https://brainchip.com/wp-content/uploads/2025/08/Akida-AKD1000-SoC-Product-Brief-V2.3-Aug.25.pdf. Accessed 2025-09-07.
- Alessio Carpegna, Alessandro Savino, and Stefano Di Carlo. Spiker+: a framework for the generation of efficient spiking neural networks fpga accelerators for inference at the edge. *IEEE Transactions on Emerging Topics in Computing*, pp. 1–15, 2024. doi: 10.1109/TETC.2024.3511676.
- Kaiwei Che, Luziwei Leng, Kaixuan Zhang, Jianguo Zhang, Max Q.-H. Meng, Jie Cheng, Qinghai Guo, and Jiangxing Liao. Differentiable hierarchical and surrogate gradient search for spiking neural networks. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc. ISBN 9781713871088.
- Ching-Yuan Chen and Krishnendu Chakrabarty. Pruning of deep neural networks for fault-tolerant memristor-based accelerators. In 2021 58th ACM/IEEE Design Automation Conference (DAC), pp. 889–894. IEEE Press, 2021. doi: 10.1109/DAC18074.2021.9586269. URL https://doi.org/10.1109/DAC18074.2021.9586269.

Lerong Chen, Jiawen Li, Yiran Chen, Qiuping Deng, Jiyuan Shen, Xiaoyao Liang, and Li Jiang. Accelerator-friendly neural-network training: Learning variations and defects in rram crossbar. In *Design, Automation & Test in Europe Conference & Exhibition (DATE)*, 2017, pp. 19–24, 2017. doi: 10.23919/DATE.2017.7926952.

Wei Fang, Yanqi Chen, Jianhao Ding, Zhaofei Yu, Timothée Masquelier, Ding Chen, Liwei Huang, Huihui Zhou, Guoqi Li, and Yonghong Tian. Spikingjelly: An open-source machine learning infrastructure platform for spike-based intelligence. *Science Advances*, 9(40):eadi1480, 2023. doi: 10.1126/sciadv.adi1480. URL https://www.science.org/doi/abs/10.1126/sciadv.adi1480.

- Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization, 2021. URL https://arxiv.org/abs/2010.01412.
- Luíza C. Garaffa, Abdullah Aljuffri, Cezar Reinbrecht, Said Hamdioui, Mottaqiallah Taouil, and Johanna Sepulveda. Revealing the secrets of spiking neural networks: The case of izhikevich neuron. In 2021 24th Euromicro Conference on Digital System Design (DSD), pp. 514–518, 2021. doi: 10.1109/DSD53832.2021.00083.
- Yufei Guo, Yuanpei Chen, Zecheng Hao, Weihang Peng, Zhou Jie, Yuhan Zhang, Xiaode Liu, and Zhe Ma. Take a shortcut back: Mitigating the gradient vanishing for training spiking neural networks. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), Advances in Neural Information Processing Systems, volume 37, pp. 24849–24867. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/2c37c5bcef24b9541550261dcd63261b-Paper-Conference.pdf.
- Zhuangyu Han, A N M Nafiul Islam, and Abhronil Sengupta. Astromorphic self-repair of neuromorphic hardware systems. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37 (6):7821–7829, Jun. 2023. doi: 10.1609/aaai.v37i6.25947. URL https://ojs.aaai.org/index.php/AAAI/article/view/25947.
- Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963. doi: 10.1080/01621459.1963.10500830.
- Hakcheon Jeong, Seungjae Han, See-On Park, Tae Ryong Kim, Jongmin Bae, Taehwan Jang, Yoonho Cho, Seokho Seo, Hyun-Jun Jeong, Seungwoo Park, Taehoon Park, Juyoung Oh, Jeongwoo Park, Kwangwon Koh, Kang-Ho Kim, Dongsuk Jeon, Inyong Kwon, Young-Gyu Yoon, and Shinhyun Choi. Self-supervised video processing with self-calibration on an analogue computing platform based on a selector-less memristor array. *Nature Electronics*, Jan 2025. ISSN 2520-1131. doi: 10.1038/s41928-024-01318-6. URL https://doi.org/10.1038/s41928-024-01318-6.
- Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-gradient methods under the polyak-łojasiewicz condition, 2020. URL https://arxiv.org/abs/1608.04636.
- James F. Kurose and Keith W. Ross. *Computer Networking: A Top-Down Approach (6th Edition)*. Pearson, 6th edition, 2012. ISBN 0132856204.
- Corey Lammie, Wei Xiang, Bernabé Linares-Barranco, and Mostafa Rahimi Azghadi. Memtorch: An open-source simulation framework for memristive deep learning systems. *Neurocomputing*, 485:124–133, 2022. ISSN 0925-2312. doi: https://doi.org/10.1016/j.neucom. 2022.02.043. URL https://www.sciencedirect.com/science/article/pii/S0925231222002053.
- Manuel Le Gallo, Corey Lammie, Julian Büchel, Fabio Carta, Omobayode Fagbohungbe, Charles Mackin, Hsinyu Tsai, Vijay Narayanan, Abu Sebastian, Kaoutar El Maghraoui, and Malte J. Rasch. Using the ibm analog in-memory hardware acceleration kit for neural network training and inference. *APL Machine Learning*, 1(4), November 2023. ISSN 2770-9019. doi: 10.1063/5.0168089. URL http://dx.doi.org/10.1063/5.0168089.

- Hyun-Jong Lee and Jae-Han Lim. Analysis on effects of fault elements in memristive neuromorphic systems. In *IJCAI 2023 GLOW workshop*, 2023. URL https://arxiv.org/abs/2312.04840.
 - Hyun-Jong Lee and Jae-Han Lim. Adaptive synaptic adjustment mechanism to improve learning performances of spiking neural networks. *Computational Intelligence*, 40(5):e70001, 2024. doi: https://doi.org/10.1111/coin.70001. URL https://onlinelibrary.wiley.com/doi/abs/10.1111/coin.70001.
 - Yudong Li, Yunlin Lei, and Xu Yang. Spikeformer: Training high-performance spiking neural network with transformer. *Neurocomputing*, 574:127279, 2024a. ISSN 0925-2312. doi: https://doi.org/10.1016/j.neucom.2024.127279. URL https://www.sciencedirect.com/science/article/pii/S092523122400050X.
 - Yuhang Li, Tamar Geller, Youngeun Kim, and Priyadarshini Panda. Seenn: towards temporal spiking early-exit neural networks. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2024b. Curran Associates Inc.
 - Shuang Lian, Jiangrong Shen, Qianhui Liu, Ziming Wang, Rui Yan, and Huajin Tang. Learnable surrogate gradient for direct training spiking neural networks. In Edith Elkind (ed.), *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI-23*, pp. 3002–3010. International Joint Conferences on Artificial Intelligence Organization, 8 2023. doi: 10.24963/ijcai.2023/335. URL https://doi.org/10.24963/ijcai.2023/335. Main Track.
 - Ling Liang, Xing Hu, Lei Deng, Yujie Wu, Guoqi Li, Yufei Ding, Peng Li, and Yuan Xie. Exploring adversarial attack in spiking neural networks with spike-compatible gradient. *IEEE Transactions on Neural Networks and Learning Systems*, 34(5):2569–2583, 2023. doi: 10.1109/TNNLS.2021. 3106961.
 - Tao Liu, Wujie Wen, Lei Jiang, Yanzhi Wang, Chengmo Yang, and Gang Quan. A fault-tolerant neural network architecture. In 2019 56th ACM/IEEE Design Automation Conference (DAC), pp. 1–6, 2019.
 - Tongliang Liu, Gábor Lugosi, Gergely Neu, and Dacheng Tao. Algorithmic stability and hypothesis complexity, 2017. URL https://arxiv.org/abs/1702.08712.
 - David G. Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, 60(2):91–110, Nov 2004. ISSN 1573-1405. doi: 10.1023/B:VISI.0000029664. 99615.94. URL https://doi.org/10.1023/B:VISI.0000029664.99615.94.
 - Abhishek Moitra, Abhiroop Bhattacharjee, Runcong Kuang, Gokul Krishnan, Yu Cao, and Priyadarshini Panda. Spikesim: An end-to-end compute-in-memory hardware evaluation tool for benchmarking spiking neural networks. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 42(11):3815–3828, 2023. doi: 10.1109/TCAD.2023.3274918.
 - Emre O. Neftci, Hesham Mostafa, and Friedemann Zenke. Surrogate gradient learning in spiking neural networks: Bringing the power of gradient-based optimization to spiking neural networks. *IEEE Signal Processing Magazine*, 36(6):51–63, 2019. doi: 10.1109/MSP.2019.2931595.
 - Yurii Nesterov. Introductory Lectures on Convex Optimization: A Basic Course. Springer Publishing Company, Incorporated, 1 edition, 2014. ISBN 1461346916.
 - Michael Pfeiffer and Thomas Pfeil. Deep learning with spiking neurons: Opportunities and challenges. *Frontiers in Neuroscience*, 12, 2018. ISSN 1662-453X. doi: 10.3389/fnins. 2018.00774. URL https://www.frontiersin.org/journals/neuroscience/articles/10.3389/fnins.2018.00774.
 - Rachmad Vidya Wicaksana Putra, Muhammad Abdullah Hanif, and Muhammad Shafique. Softsnn: low-cost fault tolerance for spiking neural network accelerators under soft errors. In *Proceedings of the 59th ACM/IEEE Design Automation Conference*, DAC '22, pp. 151–156, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450391429. doi: 10.1145/3489517. 3530657. URL https://doi.org/10.1145/3489517.3530657.

- Rachmad Vidya Wicaksana Putra, Muhammad Abdullah Hanif, and Muhammad Shafique. Rescuesnn: enabling reliable executions on spiking neural network accelerators under permanent faults. *Frontiers in Neuroscience*, 17, 2023. ISSN 1662-453X. doi: 10.3389/fnins. 2023.1159440. URL https://www.frontiersin.org/journals/neuroscience/articles/10.3389/fnins.2023.1159440.
 - Reyes-Ortiz, Jorge, Anguita, Davide, Ghio, Alessandro, Oneto, Luca, Parra, and Xavier. Human activity recognition using smartphones. UCI Machine Learning Repository, 2013. DOI: https://doi.org/10.24432/C54S4K.
 - Jozsef Sandor. The Schur-convexity of Stolarsky and Gini means. *Banach Journal of Mathematical Analysis*, 1(2):212 215, 2007. doi: 10.15352/bjma/1240336218. URL https://doi.org/10.15352/bjma/1240336218.
 - Catherine D. Schuman, Shruti R. Kulkarni, Maryam Parsa, J. Parker Mitchell, Prasanna Date, and Bill Kay. Opportunities for neuromorphic computing algorithms and applications. *Nature Computational Science*, 2(1):10–19, Jan 2022. doi: 10.1038/s43588-021-00184-y. URL https://doi.org/10.1038/s43588-021-00184-y.
 - Sumit Bam Shrestha and Garrick Orchard. Slayer: Spike layer error reassignment in time, 2018. URL https://arxiv.org/abs/1810.08646.
 - Ayesha Siddique and Khaza Anuarul Hoque. Improving reliability of spiking neural networks through fault aware threshold voltage optimization, 2023. URL https://arxiv.org/abs/2301.05266.
 - Steven S. Skiena. *The Algorithm Design Manual*. Springer, 2 edition, 2008. doi: 10.1007/978-1-84800-070-4.
 - Theofilos Spyrou, Sarah A. El-Sayed, Engin Afacan, Luis A. Camuñas-Mesa, Bernabé Linares-Barranco, and Haralampos-G. Stratigopoulos. Neuron fault tolerance in spiking neural networks. In 2021 Design, Automation & Test in Europe Conference & Exhibition (DATE), pp. 743–748, 2021. doi: 10.23919/DATE51398.2021.9474081.
 - Shigeyuki Takano and Hideharu Amano. Reconfiguration cost for reconfigurable computing architectures. In 2022 23rd ACIS International Summer Virtual Conference on Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing (SNPD-Summer), pp. 62–67, 2022. doi: 10.1109/SNPD-Summer57817.2022.00019.
 - M. L. Varshika, Sarah Johari, Jayanth Dubey, and Anup Das. Design of a tunable astrocyte neuromorphic circuitry with adaptable fault tolerance. In 2023 IEEE 66th International Midwest Symposium on Circuits and Systems (MWSCAS), pp. 904–908, 2023. doi: 10.1109/MWSCAS57524.2023. 10405978.
 - Elena-Ioana Vatajelu, Giorgio Di Natale, and Lorena Anghel. Special session: Reliability of hardware-implemented spiking neural networks (snn). In 2019 IEEE 37th VLSI Test Symposium (VTS), pp. 1–8, 2019. doi: 10.1109/VTS.2019.8758653.
 - Yeshwanth Venkatesha, Youngeun Kim, Leandros Tassiulas, and Priyadarshini Panda. Federated learning with spiking neural networks. *IEEE Transactions on Signal Processing*, 69:6183–6194, 2021. doi: 10.1109/TSP.2021.3121632.
 - The H. Vu, Ogbodo Mark Ikechukwu, and Abderazek Ben Abdallah. Fault-tolerant spike routing algorithm and architecture for three dimensional noc-based neuromorphic systems. *IEEE Access*, 7:90436–90452, 2019. doi: 10.1109/ACCESS.2019.2925085.
 - Ziming Wang, Runhao Jiang, Shuang Lian, Rui Yan, and Huajin Tang. Adaptive smoothing gradient learning for spiking neural networks. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 35798–35816. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/wang23j.html.

- Rachmad Vidya Wicaksana Putra, Muhammad Abdullah Hanif, and Muhammad Shafique. Respawn: Energy-efficient fault-tolerance for spiking neural networks considering unreliable memories. In 2021 IEEE/ACM International Conference On Computer Aided Design (ICCAD), pp. 1–9, 2021. doi: 10.1109/ICCAD51958.2021.9643524.
- Torbjörn Wigren and Reem Karaki. Globally stable wireless data flow control. *IEEE Transactions on Control of Network Systems*, 5(1):469–478, 2018. doi: 10.1109/TCNS.2016.2619906.
- Shuangming Yang, Jiang Wang, Bin Deng, Mostafa Rahimi Azghadi, and Bernabe Linares-Barranco. Neuromorphic context-dependent learning framework with fault-tolerant spike routing. *IEEE Transactions on Neural Networks and Learning Systems*, 33(12):7126–7140, 2022. doi: 10.1109/TNNLS.2021.3084250.
- Man Yao, Guangshe Zhao, Hengyu Zhang, Yifan Hu, Lei Deng, Yonghong Tian, Bo Xu, and Guoqi Li. Attention spiking neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(8):9393–9410, 2023. doi: 10.1109/TPAMI.2023.3241201.
- Shlomo Yitzhaki and Edna Schechtman. *The Gini Methodology: A Primer on a Statistical Methodology*. Springer Series in Statistics. Springer, New York, 2013. ISBN 978-1-4614-4719-1. doi: 10.1007/978-1-4614-4720-7.
- Anlan Yu, Ning Lyu, Jieming Yin, Zhiyuan Yan, and Wujie Wen. COLA: Orchestrating error coding and learning for robust neural network inference against hardware defects. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 40277–40289. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/yu23a.html.
- Aybars Yunusoglu, Dexter Le, Murat Isik, I. Can Dikmen, and Teoman Karadag. Neuromorphic circuits with spiking astrocytes for increased energy efficiency, fault tolerance, and memory capacitance, 2025. URL https://arxiv.org/abs/2502.20492.
- Friedemann Zenke and Tim P. Vogels. The remarkable robustness of surrogate gradient learning for instilling complex function in spiking neural networks. *Neural Computation*, 33(4):899–925, 03 2021. ISSN 0899-7667. doi: 10.1162/neco_a_01367. URL https://doi.org/10.1162/neco_a_01367.

A ADDITIONAL EXPERIMENTAL RESULTS ON CLASSIFICATION ACCURACY

We present additional experimental results that support the proposed mechanism in this section. The additional results demonstrate that our mechanism enhances the neural networks' fault tolerance more effectively than the existing methods in various environments and scenarios, including real hardware. Furthermore, we discuss the changes in the proposed mechanism's fault mitigation ability by adopting various settings to our mechanism.

A.1 ADDITIONAL DATASETS BEYOND MNIST, FMNIST, CIFAR-10, AND CIFAR-100

We use UCI-HAR and Tiny-ImageNet to evaluate the fault mitigation ability of our mechanism on a sequential and large-scale dataset.

A.1.1 SEQUENTIAL DATASET

To demonstrate that our mechanism works well with the models using sequential datasets, we conduct experiments with UCI-HAR, which comprises six types of human activities collected by electric sensors, and a dataset consisting of verbal sounds of digits (Reyes-Ortiz et al., 2013).

Table 4: The small MLP models' classification accuracy in a 95% confidence interval using UCI-HAR under SAFs with 2 time steps.

FAULT RATIO(%)	BASELINE	ECOC	SOFTSNN	ROUTING	ASTROCYTE	FALVOLT	LIFA	PROPOSED
0	64.14 ± 4.91	69.86 ± 4.73	63.41 ± 4.78	65.17 ± 4.65	64.78 ± 4.39	65.76 ± 4.55	64.69 ± 4.81	63.14 ± 4.8
10	20.19 ± 2.95	23.01 ± 4.28	22.21 ± 4.63	23.1 ± 4.58	21.25 ± 3.56	23.24 ± 4.01	16.93 ± 4.52	50.31 ± 4.68
20	19.66 ± 2.54	17.1 ± 0	20.52 ± 3.85	17.1 ± 0	20.59 ± 2.81	22.34 ± 3.68	17.1 ± 0	$\textbf{49.24} \pm \textbf{4.59}$
30	17.1 ± 0	17.1 ± 0	19.9 ± 2.26	17.1 ± 0	17.1 ± 0	19.85 ± 2.92	17.1 ± 0	$\textbf{47.31} \pm \textbf{4.36}$
40	17.1 ± 0	16.93 ± 0	17.1 ± 0	$\textbf{17.21} \pm \textbf{0.09}$				
50	17.1 ± 0	16.93 ± 0	17.1 ± 0	17.1 ± 0				
60	17.1 ± 0	16.93 ± 0	17.1 ± 0	17.1 ± 0				
70	17.1 ± 0	16.93 ± 0	17.1 ± 0	17.1 ± 0				
80	17.1 ± 0	16.93 ± 0	17.1 ± 0	17.1 ± 0				
90	17.1 ± 0	17.1 ± 0	17.1 ± 0	17.1 ± 0	17.1 ± 0	16.93 ± 0	17.1 ± 0	17.1 ± 0

Table 4 shows the classification accuracy of the MLP model using UCI-HAR. Using the MLP model with a sequential dataset, the model with the proposed mechanism exhibits better fault tolerance than the baseline and benchmarks, classifying data samples more accurately than the models with the baseline and benchmarks. This result shows that our mechanism works well with the sequential dataset.

A.1.2 LARGE IMAGE DATASET

We use the ResNet-34 model to classify data samples in Tiny-ImageNet, which is a small version of the ImageNet dataset, consisting of 64×64 pixel images with 200 classes. We measure the classification accuracy of the model with Tiny-ImageNet under SAFs.

Table 5: The ResNet-34 models' classification accuracy in a 95% confidence interval using Tiny-ImageNet under SAFs with 2 time steps.

FAULT RATIO(%)	BASELINE	ECOC	SOFTSNN	ROUTING	ASTROCYTE	FALVOLT	LIFA	PROPOSED
0	3.38 ± 2.11	3.61 ± 2.24	3.27 ± 2.08	3.24 ± 2.16	0.65 ± 0.15	3.47 ± 1.98	0.59 ± 0.09	3.96 ± 2.39
10	0.65 ± 0.12	0.69 ± 0.15	0.68 ± 0.13	0.71 ± 0.19	0.5 ± 0	0.72 ± 0.16	0.5 ± 0	1.17 ± 0.28
20	0.5 ± 0	$\textbf{0.63} \pm \textbf{0.12}$						
30	0.5 ± 0	$\textbf{0.56} \pm \textbf{0.06}$						
40	0.5 ± 0							
50	0.5 ± 0							
60	0.5 ± 0							
70	0.5 ± 0							
80	0.5 ± 0							
90	0.5 ± 0							

Table 5 presents the average classification accuracy of the ResNet-34 model with the baseline, benchmarks, and proposed mechanism using Tiny-ImageNet. The classification accuracy of the model degrades because the dataset is complex, and SNN models have lower learning capabilities compared to DNN models. Despite the low classification accuracy, the model with our mechanism classifies data samples in Tiny-ImageNet with the highest accuracy. Moreover, the model with our mechanism exhibits higher accuracy than others in the clean scenario (without SAFs). This is because our mechanism leads the models to emit the output precisely through entropy-based decoding.

A.2 CHANGING THE NUMBER OF TIME STEPS

We change the number of time steps to 4 and 8, observing the accuracy trend of all SNN models in the number of time steps. We obtain the experiment results with 4 and 8 time steps by repeating the experiments 10 times.

A.2.1 4 TIME STEPS

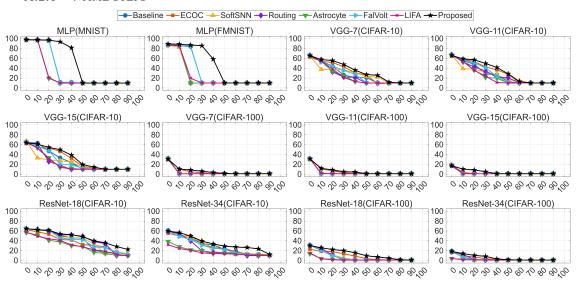


Figure 6: Average classification accuracy of various SNN models with the baseline, benchmarks, and proposed mechanism under SAFs using 4 time steps. The x-axis is the fault ratio (%) and the y-axis is the accuracy (%).

Figure 6 illustrates the classification accuracy of SNN models compared to the baseline, benchmarks, and proposed mechanism under SAFs with 4 time steps. The models with our mechanism exhibit the best fault tolerance to faults in most cases, like the experimental results with 2 time steps. We observe that the classification accuracy of all models overall improves as the number of time steps increases, because the large number of time steps improves the performance of SNNs (Li et al., 2024b).

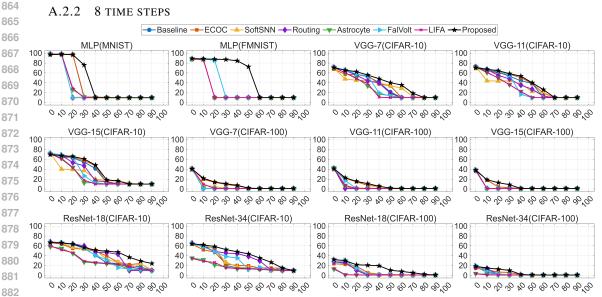


Figure 7: Average classification accuracy of various SNN models with the baseline, benchmarks, and proposed mechanism under SAFs using 8 time steps. The x-axis is the fault ratio (%) and the y-axis is the accuracy (%).

Figure 7 depicts the classification accuracy of SNN models compared to the baseline, benchmarks, and proposed mechanism under SAFs with 8 time steps. As demonstrated in the experimental results with 2 and 4 time steps, the models with our mechanism classify data samples most accurately. The models' accuracy is also higher than when using 2 and 4 time steps. Interestingly, the fault mitigation ability of our mechanism degrades only in the experiment with MNIST samples. This is because the MNIST samples contain fewer pixels than FMNIST, CIFAR-10, and CIFAR-100. However, they are divided into too small fragments, and these fragments do not have sufficient information for the MLP model to learn. Thus, the fault tolerance of the model with our mechanism weakens, although it is more fault-robust than models with the benchmarks.

A.3 Under the different types of synaptic faults

We inject RWFs and CEFs into the synapses of the SNN models and measure the fault mitigation ability of the benchmarks and the proposed mechanism. The models with our mechanism classify data samples most accurately under RWFs and CEFs.

A.3.1 RWFs

We use a Gaussian distribution to model RWFs, setting the standard deviation of the distribution to 0.5 (Garaffa et al., 2021; Spyrou et al., 2021; Vatajelu et al., 2019). We

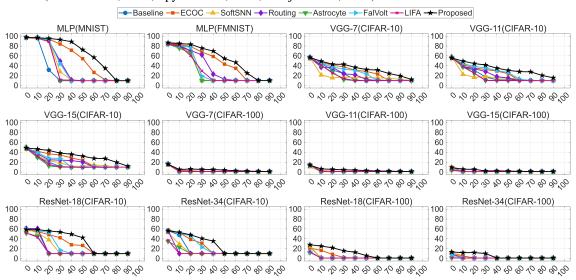


Figure 8: Average classification accuracy of various SNN models with the baseline, benchmarks, and proposed mechanism under RWFs using 2 time steps. The x-axis is the fault ratio (%) and the y-axis is the accuracy (%).

Figure 8 depicts the classification accuracy of SNN models with the baseline, benchmarks, and proposed mechanism under RWFs. The models with our mechanism exhibit the highest accuracy in classifying MNIST, FMNIST, CIFAR-10, and CIFAR-100. This is because our mechanism successfully prevents the pre-activation magnitude from increasing excessively by RWFs. Interestingly, ECOC presents high fault mitigation ability under RWFs. This is because ECOC uses error correcting codes, which are robust against Gaussian noise in channels to compensate for errors caused by faults in the last layer (Liu et al., 2019).

A.3.2 CEFs

 CEFs change the connections between spiking neurons randomly, ruining the learned information of SNN models (Vatajelu et al., 2019).

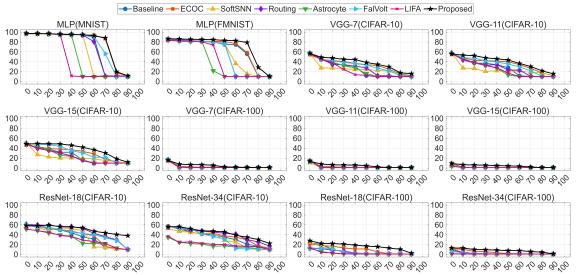


Figure 9: Average classification accuracy of various SNN models with the baseline, benchmarks, and proposed mechanism under CEFs using 2 time steps. The x-axis is the fault ratio (%) and the y-axis is the accuracy (%).

Figure 9 illustrates the classification accuracy of SNN models with the baseline, benchmarks, and proposed mechanism under CEFs. Our mechanism also presents the best fault mitigation ability. The classification accuracy of SNN models under CEFs is higher than that of the models under SAFs and RWFs. This is because the weights of faulty synapses are uniform under CEFs, and the pre-activation magnitude does not increase significantly. Thus, the pre-activation does not lie in a value that is far from the surrogate gradient corridor.

A.4 ABLATION STUDY ON THE COMBINATION OF OUR MECHANISM

We conduct ablation studies by changing the settings of our mechanism (horizontally-fixed vs Ginibased and only complexity-based sensitivity score vs complexity and influence combined sensitivity score).

A.4.1 MLP MODEL (MNIST, FMNIST)

Table 6: The MLP model's classification accuracy in a 95% confidence interval with different settings of our mechanism under SAFs.

FAULT RATIO(%)	BASELINE	HORIZONTAL	GINI(COMPLEXITY ONLY)	GINI(PROPOSED)
	Acc	CURACY (%) WIT	H MLP (MNIST)	
0	97.35 ± 0.46	86.68 ± 3.51	97.37 ± 0.44	97.44 ± 0.39
10	96.94 ± 1.01	75.77 ± 4.96	96.76 ± 0.95	$\textbf{97.25} \pm \textbf{0.98}$
20	11.35 ± 0	71.13 ± 4.81	93.64 ± 1.44	$\textbf{93.84} \pm \textbf{1.37}$
30	10.82 ± 0	65.78 ± 5.39	91.56 ± 2.61	$\textbf{93.55} \pm \textbf{1.94}$
40	10.82 ± 0	11.35 ± 0	11.35 ± 0	$\textbf{86.71} \pm \textbf{4.68}$
50	9.8 ± 0	9.8 ± 0	9.8 ± 0	$\textbf{26.91} \pm \textbf{7.53}$
60	9.8 ± 0	9.8 ± 0	9.8 ± 0	9.8 ± 0
70	9.8 ± 0	9.8 ± 0	9.8 ± 0	9.8 ± 0
80	9.8 ± 0	9.8 ± 0	9.8 ± 0	9.8 ± 0
90	9.8 ± 0	9.8 ± 0	9.8 ± 0	9.8 ± 0
	Acc	URACY (%) WITH	H MLP (FMNIST)	
0	83.99 ± 0.86	78.09 ± 1.03	86.29 ± 0.94	$\textbf{86.8} \pm \textbf{0.89}$
10	83.74 ± 1.79	76.78 ± 1.9	85.49 ± 1.43	$\textbf{85.6} \pm \textbf{1.23}$
20	10 ± 0	74.22 ± 2.54	80.41 ± 2.19	$\textbf{85.33} \pm \textbf{1.45}$
30	10 ± 0	70.65 ± 6.17	79.6 ± 3.72	83.53 ± 2.69
40	10 ± 0	17.16 ± 8.84	21.93 ± 7.63	$\textbf{54.55} \pm \textbf{6.76}$
50	10 ± 0	10 ± 0	10 ± 0	10 ± 0
60	10 ± 0	10 ± 0	10 ± 0	10 ± 0
70	10 ± 0	10 ± 0	10 ± 0	10 ± 0
80	10 ± 0	10 ± 0	10 ± 0	10 ± 0
90	10 ± 0	10 ± 0	10 ± 0	10 ± 0

Table 6 presents the classification accuracy of the baseline, horizontally-fixed fragmentation, Ginibased fragmentation using image complexity for the sensitivity score, and Gini-based fragmentation using image complexity and fault influence of the first layer for the sensitivity score (proposed) when using MNIST models to classify MNIST and FMNIST data samples. The proposed version significantly enhances the fault tolerance of the MLP model, as demonstrated by its performance on MNIST and FMNIST, compared to other settings. This is because the MLP model is vulnerable to faults in the first layer, as mentioned in Section 4. Thus, using fault influence for the sensitivity score enhances our mechanism's fault mitigation ability since it induces the mechanism to minimize the pre-activation magnitude. We also observe that the wrong fragmentation strategy degrades classification performance because it damages the information of data samples.

A.4.2 VGG-7 AND RESNET-18 MODELS (CIFAR-10 AND CIFAR-100)

Table 7: The VGG-7 and ResNet-18 models' classification accuracy in a 95% confidence interval with different settings of our mechanism under SAFs.

FAULT RATIO(%)	BASELINE	HORIZONTAL	GINI(COMPLEXITY ONLY)	GINI(PROPOSED)
	Accur	RACY (%) WITH	VGG-7 (CIFAR-10)	
0	56.26 ± 1.28	55.75 ± 1.33	56.68 ± 1.28	$\textbf{56.86} \pm \textbf{1.63}$
10	40.09 ± 3.34	48.8 ± 2.67	49.13 ± 2.32	$\textbf{50.38} \pm \textbf{2.07}$
20	31.78 ± 4.97	42.69 ± 3.83	45.39 ± 3.44	$\textbf{45.58} \pm \textbf{3.76}$
30	11.74 ± 4.81	36.55 ± 4.64	38.23 ± 4.91	$\textbf{40.26} \pm \textbf{5.04}$
40	10.79 ± 3.56	29.09 ± 5.72	31.18 ± 5.52	$\textbf{32.07} \pm \textbf{5.15}$
50	10 ± 0	10 ± 0	10 ± 0	$\textbf{17.08} \pm \textbf{5.74}$
60	10 ± 0	10 ± 0	10 ± 0	10 ± 0
70	10 ± 0	10 ± 0	10 ± 0	10 ± 0
80	10 ± 0	10 ± 0	10 ± 0	10 ± 0
90	10 ± 0	10 ± 0	10 ± 0	10 ± 0
	Accurac	CY (%) WITH RE	SNET-18 (CIFAR-100)	
0	23.77 ± 4.56	24.12 ± 4.69	25.11 ± 4.74	27.96 ± 4.83
10	10.64 ± 3.82	23.98 ± 4.41	24.15 ± 4.32	$\textbf{25.5} \pm \textbf{4.79}$
20	6.1 ± 2.65	19.86 ± 5.52	21.92 ± 5.48	$\textbf{24.41} \pm \textbf{5.38}$
30	1.98 ± 0.93	16.62 ± 4.63	18.43 ± 5.01	$\textbf{19.72} \pm \textbf{5.23}$
40	1.03 ± 0.82	13.53 ± 4.29	14.26 ± 4.57	$\textbf{16.8} \pm \textbf{4.72}$
50	0.97 ± 0.79	11.89 ± 4.17	12.34 ± 4.48	$\textbf{13.13} \pm \textbf{4.52}$
60	1.01 ± 0.86	4.68 ± 2.1	5.12 ± 2.96	$\textbf{6.25} \pm \textbf{3.26}$
70	1 ± 0	2.25 ± 1.12	3.9 ± 1.44	$\textbf{4.27} \pm \textbf{1.56}$
80	1 ± 0	1 ± 0	1 ± 0	$\textbf{1.24} \pm \textbf{0.81}$
90	1 ± 0	1 ± 0	1 ± 0	1 ± 0

Table 7 exhibits our ablation study with VGG-7 and ResNet-18 models using CIFAR-10 and CIFAR-100. Although the proposed mechanism outperforms other settings, the improvement in fault mitigation ability is not as large as in the cases with the MLP models. This is because VGG-7 and ResNet-18 models are not as vulnerable to faults in the first layer as the MLP model since they have additional features to compensate for errors during gradient calculations. Therefore, containing the fault influence in the sensitivity score does not significantly enhance the fault mitigation ability of our mechanism.

A.5 VARIOUS HYPERPARAMETER SETTINGS WITH OUR MECHANISM

We change (λ_s, λ_w) to (0.5, 0.5), (0.5, 0.1), and (0.1, 0.5) to show the influence of hyperparameters on our mechanism.

A.5.1 MLP MODEL (MNIST AND FMNIST)

Table 8: The MLP models' classification accuracy in a 95% confidence interval with different settings of the hyperparameter (λ_s, λ_w) under SAFs.

FAULT RATIO(%)	(0.5, 0.5)	(0.5, 0.1)	(0.1, 0.5)	(0.1, 0.1) (DEFAULT)
	Accura	CY (%) WITH M	LP (MNIST)	
0	97.28 ± 0.53	97.35 ± 0.54	97.29 ± 0.51	97.44 ± 0.39
10	96.34 ± 1.25	96.58 ± 1.38	96.05 ± 1.19	$\textbf{97.25} \pm \textbf{0.98}$
20	95.29 ± 1.67	93.3 ± 1.63	95.43 ± 1.58	$\textbf{95.84} \pm \textbf{1.37}$
30	91.59 ± 2.82	92.15 ± 2.97	86.93 ± 2.86	93.55 ± 1.94
40	9.8 ± 0	14.28 ± 5.47	10.48 ± 1.39	$\textbf{86.71} \pm \textbf{4.68}$
50	9.8 ± 0	12.6 ± 5.61	9.8 ± 0	$\textbf{26.91} \pm \textbf{7.53}$
60	9.8 ± 0	9.8 ± 0	9.8 ± 0	9.8 ± 0
70	9.8 ± 0	9.8 ± 0	9.8 ± 0	9.8 ± 0
80	9.8 ± 0	9.8 ± 0	9.8 ± 0	9.8 ± 0
90	9.8 ± 0	9.8 ± 0	9.8 ± 0	9.8 ± 0
	Accurac	CY (%) WITH MI	LP (FMNIST)	
0	85.53 ± 1.05	85.56 ± 1.14	85.29 ± 0.99	$\textbf{86.8} \pm \textbf{0.89}$
10	85.46 ± 1.46	85.52 ± 1.31	85.23 ± 1.36	$\textbf{85.6} \pm \textbf{1.23}$
20	80.89 ± 3.48	83.57 ± 2.12	84.41 ± 2.29	$\textbf{85.33} \pm \textbf{1.45}$
30	80.23 ± 3.8	31.58 ± 7.03	75.39 ± 5.45	83.53 ± 2.69
40	33.54 ± 6.15	18.84 ± 7.65	10 ± 0	$\textbf{54.55} \pm \textbf{6.76}$
50	10 ± 0	10 ± 0	10 ± 0	10 ± 0
60	10 ± 0	10 ± 0	10 ± 0	10 ± 0
70	10 ± 0	10 ± 0	10 ± 0	10 ± 0
80	10 ± 0	10 ± 0	10 ± 0	10 ± 0
90	10 ± 0	10 ± 0	10 ± 0	10 ± 0

Table 8 exhibits the MLP model's classification accuracy by changing (λ_s and λ_w) under SAFs. The hyperparameter settings predominantly affect the classification accuracy of the MLP model, as they adjust how the mechanism mitigates the adverse influence of synaptic faults in the first layer, which damages the model the most severely. We observe that increasing the hyperparameters and strengthening the effects of the fault influence map do not always leverage the MLP model's fault tolerance. This is because the excessive effects of the fault influence map prevent our mechanism from setting the angle accurately by reflecting the complexity of the input samples and the fault influence in a balanced way. In addition, we observe that the weight projection map affects our mechanism more predominantly than the saliency map, since the weight projection map is more sensitive to changes in weights due to synaptic faults.

A.5.2 VGG-7 AND RESNET-18 MODELS (CIFAR-10 AND CIFAR-100)

Table 9: The VGG-7 and ResNet-18 models' classification accuracy in a 95% confidence interval with different settings of the hyperparameter (λ_s , λ_w) under SAFs.

FAULT RATIO(%)	(0.5, 0.5)	(0.5, 0.1)	(0.1, 0.5)	(0.1, 0.1) (DEFAULT)
The Er Kille (70)				
	ACCURACY	(%) WITH VGC	G-7 (CIFAR-10)	
0	56.48 ± 1.84	56.06 ± 1.73	56.03 ± 1.89	$\textbf{56.86} \pm \textbf{1.63}$
10	48.12 ± 3.04	48.88 ± 2.85	49.34 ± 2.16	$\textbf{50.38} \pm \textbf{2.07}$
20	44.68 ± 3.69	44.56 ± 3.66	45.3 ± 3.59	$\textbf{45.58} \pm \textbf{3.76}$
30	39.42 ± 5.25	39.22 ± 5.59	39.92 ± 5.23	$\textbf{40.26} \pm \textbf{5.04}$
40	29.62 ± 5.38	29.52 ± 5.04	30.21 ± 5.4	$\textbf{32.07} \pm \textbf{5.15}$
50	15.09 ± 4.76	15.28 ± 5.17	15.31 ± 4.98	$\textbf{17.08} \pm \textbf{5.74}$
60	10 ± 0	10 ± 0	10 ± 0	10 ± 0
70	10 ± 0	10 ± 0	10 ± 0	10 ± 0
80	10 ± 0	10 ± 0	10 ± 0	10 ± 0
90	10 ± 0	10 ± 0	10 ± 0	10 ± 0
	ACCURACY (%	%) WITH RESNE	T-18 (CIFAR-10	00)
0	26.96 ± 5.08	26.02 ± 4.93	25.71 ± 5.86	27.96 ± 4.83
10	24.26 ± 4.29	24.76 ± 4.38	21.75 ± 4.71	$\textbf{25.5} \pm \textbf{4.79}$
20	20.09 ± 5.66	21.23 ± 5.96	19.03 ± 5.09	$\textbf{24.41} \pm \textbf{5.38}$
30	16.88 ± 5.48	17.7 ± 5.48	15.89 ± 5.42	$\textbf{19.72} \pm \textbf{5.23}$
40	13.25 ± 5.24	13.68 ± 5.05	13.25 ± 4.93	$\textbf{16.8} \pm \textbf{4.71}$
50	12.04 ± 4.93	12.55 ± 4.34	9.14 ± 4.69	$\textbf{13.13} \pm \textbf{4.5}$
60	4.46 ± 3.52	4.77 ± 3.15	4.86 ± 2.96	$\textbf{6.25} \pm \textbf{3.26}$
70	3.61 ± 1.95	2.89 ± 1.7	3.05 ± 1.78	$\textbf{4.27} \pm \textbf{1.56}$
80	1 ± 0	1 ± 0	1 ± 0	$\textbf{1.24} \pm \textbf{0.81}$
90	1 ± 0	1 ± 0	1 ± 0	1 ± 0

Table 9 presents the VGG and ResNet Models' classification accuracy by changing (λ_s and λ_w) under SAFs. Since the significance of the faults in the first layer of the VGG and ResNet Models is weaker than that of the MLP model, the sensitivity to hyperparameter setting is smaller than that of the MLP model. Despite the low significance of tuning hyperparameters when using the VGG and ResNet Models, setting the hyperparameters to the proper value is still important to ensure the best fault mitigation ability of our mechanism.

A.6 USING DEEP NEURAL NETWORKS

We inject SAFs into synapses of Deep Neural Networks (DNNs) version of the SNN models, measuring their classification accuracy with the baseline, benchmarks, and proposed mechanism. We use Cross-Entropy (CE) as a loss function and Rectified Linear Unit (ReLU) as an activation function. We set the range of weights to [-100, 100] for MLP and [-500, 500] for CNNs (VGG and ResNet) since current DNN accelerator devices have a large weight range (Liu et al., 2019; Chen et al., 2017). Other settings are the same as the SNN models. We exclude Astrocyte, FalVolt, and LIFA from the benchmarks since they necessarily require bio-plausible spiking neuron models for operation. We use 2 fragments for our mechanism.

Table 10: The DNN models' classification accuracy in a 95% confidence interval using MNIST, FMNIST, CIFAR-10, and CIFAR-100 under SAFs.

FAULT RATIO(%)	BASELINE	ECOC	SOFTSNN (TUNED FOR DNN)	ROUTING	PROPOSED				
ACCURACY(%) WITH MLP (MNIST)									
0	98.48 ± 0.19	$\textbf{98.6} \pm \textbf{0.26}$	98.49 ± 0.23	98.45 ± 0.31	98.54 ± 0.25				
10	97.32 ± 0.86	97.28 ± 0.81	97.48 ± 0.96	97.28 ± 0.89	$\textbf{97.56} \pm \textbf{0.74}$				
20	96.62 ± 1.55	9.8 ± 0	97.02 ± 1.32	96.76 ± 1.47	$\textbf{96.85} \pm \textbf{1.38}$				
30	96.58 ± 1.93	9.8 ± 0	96.3 ± 1.89	96.64 ± 1.69	$\textbf{96.72} \pm \textbf{1.71}$				
40	96.35 ± 1.81	9.8 ± 0	96.02 ± 1.76	96.43 ± 1.74	$\textbf{96.48} \pm \textbf{1.85}$				
50	93.3 ± 3.52	9.8 ± 0	93.47 ± 3.62	93.54 ± 4.02	93.58 ± 3.91				
60	77.36 ± 9.15	9.8 ± 0	78.03 ± 8.99	79.89 ± 8.56	$\textbf{82.62} \pm \textbf{8.8}$				
70	65.9 ± 10.49	9.8 ± 0	68.39 ± 10.05	67.79 ± 9.95	$\textbf{74.8} \pm \textbf{9.72}$				
80	39.64 ± 8.77	9.8 ± 0	41.31 ± 9.12	18.71 ± 8.89	$\textbf{55.43} \pm \textbf{8.73}$				
90	8.92 ± 0	9.8 ± 0	9.8 ± 0	15.5 ± 5.7	$\textbf{19.07} \pm \textbf{6.21}$				
		ACCURACY	(%) WITH MLP (FMNIST)						
0	90.05 ± 1.22	$\textbf{91.24} \pm \textbf{1.3}$	90.08 ± 1.24	90.12 ± 1.13	90.51 ± 1.19				
10	84.38 ± 2.73	10 ± 0	84.72 ± 2.86	86.03 ± 2.69	$\textbf{86.83} \pm \textbf{2.71}$				
20	71.67 ± 3.49	10 ± 0	72.07 ± 3.35	74.51 ± 3.48	$\textbf{79.66} \pm \textbf{3.48}$				
30	68.01 ± 5.05	10 ± 0	70.75 ± 4.94	73.54 ± 4.75	$\textbf{74.63} \pm \textbf{4.79}$				
40	64.93 ± 5.91	10 ± 0	63.5 ± 6.04	64.06 ± 5.97	$\textbf{65.23} \pm \textbf{6.25}$				
50	58.63 ± 6.27	10 ± 0	62.78 ± 6.38	60.92 ± 6.09	64.8 \pm 6.31				
60	53.57 ± 7.86	10 ± 0	57.01 ± 8.11	54.51 ± 9.23	$\textbf{58.68} \pm \textbf{8.77}$				
70	38.01 ± 9.91	10 ± 0	39.48 ± 10.26	35.12 ± 9.84	$\textbf{41.62} \pm \textbf{10.34}$				
80	23.5 ± 8.44	10 ± 0	26.45 ± 8.79	25.41 ± 7.98	$\textbf{33.21} \pm \textbf{8.59}$				
90	10 ± 0	10 ± 0	10 ± 0	10 ± 0	$\textbf{11.02} \pm \textbf{1.02}$				
		ACCURACY(%) WITH VGG-7 (CIFAR-10)						
0	83.21 ± 2.76	83.36 ± 2.58	83.58 ± 2.53	83.73 ± 2.47	84.29 ± 2.63				
0.01	10 ± 0	10 ± 0	52.97 ± 4.84	67.78 ± 5.04	$\textbf{70.36} \pm \textbf{4.81}$				
0.025	10 ± 0	10 ± 0	10 ± 0	10 ± 0	10 ± 0				
0.05	10 ± 0	10 ± 0	10 ± 0	10 ± 0	10 ± 0				
0.075	10 ± 0	10 ± 0	10 ± 0	10 ± 0	10 ± 0				
		ACCURACY(%)	WITH RESNET-18 (CIFAR-100)						
0	53.09 ± 1.02	40.96 ± 1.56	53.51 ± 1.16	53.6 ± 1.24	$\textbf{53.82} \pm \textbf{1.19}$				
0.01	46.17 ± 1.35	37.87 ± 1.53	47.26 ± 1.58	46.99 ± 1.56	$\textbf{48.49} \pm \textbf{1.47}$				
0.025	43.32 ± 1.68	34.05 ± 1.79	43.76 ± 1.73	44.07 ± 1.79	$\textbf{45.09} \pm \textbf{1.66}$				
0.05	41.14 ± 1.54	32.28 ± 2.01	42.23 ± 1.69	42.51 ± 1.55	$\textbf{42.78} \pm \textbf{1.48}$				
0.075	39.05 ± 1.82	31.64 ± 1.98	40.64 ± 1.91	39.89 ± 2.05	$\textbf{41.32} \pm \textbf{1.85}$				

Table 10 presents the classification accuracy of DNN models with the baseline, benchmarks, and proposed mechanism. The models with ours exhibit the highest fault robustness among the DNN models since our mechanism prevents the pre-activation from increasing excessively, and gradients do not explode severely during training. The MLP DNN model presents higher fault tolerance than the SNN model regardless of the datasets. This is because the gradient vanishing caused by pre-activation magnitude growth does not occur severely, since gradients are active when the pre-activation is larger than 0 in ReLU. Conversely, the fault tolerance of the complicated CNN models (VGG and ResNet) degrades seriously. This situation appears because these models use lots of batch normalization layers. The normalization layers normalize the whole channels in the same scale calculated with the batch samples, causing the amplification of inputs that enter faulty synapses. However, the SNN VGG and ResNet models are more tolerant of faults than the DNN versions. This is because the spiking neurons block the perturbation from faults through their internal activation mechanism: only firing and emitting spikes when their membrane potential reaches the threshold (Liang et al., 2023).

A.7 EVALUATIONS WITH REAL FPGA HARDWARE

 We implement the MLP SNN model on a real FPGA device (AMD Virtex UltraScale+ HBM VU47P of Amazon F2 instance) with SpikerPlus, which is a powerful library to convert Python scripts for SNNs to VHSIC Hardware Description Language (VHDL) (Carpegna et al., 2024). Due to circuit-level limitations, current neuromorphic devices do not support on-chip training. Thus, we train the fault-injected SNN model in a software environment with a Graphics Processing Unit (GPU), saving the trained weights, and convert the Python script of the software-based models to the VHDL script. Then, we synthesize the FPGA circuit with Xilinx Vivado and Amazon FPGA Image (AFI), which is widely used for handling FPGAs on Amazon F2 instances. We set the floating-point for the membrane potential (neuron in the SpikerPlus paper) and bit width for the synaptic weights of the FPGA device to 16 and 8/6, referring to the settings of real hardware (BrainChipInc, 2025). Other settings are the same as the settings in Subsection 6.1.

Table 11: The FPGA-based MLP SNNs' classification accuracy in a 95% confidence interval using MNIST, FMNIST, and UCI-HAR with 2 time steps under SAFs.

FAULT RATIO(%)	BASELINE	ECOC	SOFTSNN	ROUTING	ASTROCYTE	FALVOLT	LIFA	PROPOSED	
HARDWARE-BASED MLP'S ACCURACY (%) WITH MNIST									
0	94.36 ± 0.52	93.86 ± 0.66	94.56 ± 0.63	94.34 ± 0.71	94.06 ± 0.68	93.92 ± 0.79	93.98 ± 0.85	94.59 ± 0.65	
10	93.18 ± 1.29	90.94 ± 2.35	93.29 ± 1.56	93.58 ± 1.9	93.26 ± 1.44	92.69 ± 1.62	92.59 ± 1.74	$\textbf{94.01} \pm \textbf{1.41}$	
20	11.35 ± 0	88.11 ± 2.74	9.8 ± 0	9.8 ± 0	10.1 ± 1.05	89.97 ± 2.69	15.92 ± 5.31	$\textbf{91.14} \pm \textbf{2.03}$	
30	10.82 ± 0	11.35 ± 0	9.8 ± 0	9.8 ± 0	10.82 ± 0	9.8 ± 0	9.8 ± 0	$\textbf{90.57} \pm \textbf{3.16}$	
40	9.8 ± 0	9.8 ± 0	9.8 ± 0	11.35 ± 0	9.8 ± 0	9.8 ± 0	9.8 ± 0	$\textbf{80.49} \pm \textbf{6.65}$	
50	9.8 ± 0	9.8 ± 0	9.8 ± 0	11.35 ± 0	9.8 ± 0	9.8 ± 0	9.8 ± 0	$\textbf{21.46} \pm \textbf{4.72}$	
60	9.8 ± 0								
70	9.8 ± 0								
80	9.8 ± 0								
90	9.8 ± 0								
		HARD	WARE-BASED M	LP's Accurac	Y (%) WITH FM	NIST			
0	84.21 ± 1.79	83.58 ± 1.91	83.65 ± 2.13	84.19 ± 1.82	83.9 ± 1.75	83.97 ± 1.68	84.1 ± 1.58	$\textbf{87.01} \pm \textbf{1.45}$	
10	82.16 ± 2.28	82.75 ± 2.46	80.74 ± 2.67	83.58 ± 2.33	82.86 ± 2.09	83.11 ± 2.53	82.46 ± 2.38	$\textbf{86.14} \pm \textbf{2.84}$	
20	79.6 ± 3.05	10 ± 0	10 ± 0	10 ± 0	10 ± 0	82.54 ± 3.18	17.26 ± 5.98	$\textbf{83.98} \pm \textbf{3.27}$	
30	10 ± 0	$\textbf{81.71} \pm \textbf{3.7}$							
40	10 ± 0	$\textbf{52.75} \pm \textbf{4.89}$							
50	10 ± 0								
60	10 ± 0								
70	10 ± 0								
80	10 ± 0								
90	10 ± 0								
		HARDV	VARE-BASED MI	LP's ACCURACY	(%) WITH UCI	-HAR			
0	60.14 ± 3.87	65.27 ± 3.56	60.48 ± 3.81	61.15 ± 3.69	61.78 ± 3.48	61.83 ± 3.45	61.52 ± 3.62	60.08 ± 3.24	
10	17.1 ± 0	17.1 ± 0	15.21 ± 0	17.1 ± 0	15.21 ± 0	18.34 ± 0	16.93 ± 0	47.01 ± 5.17	
20	18.52 ± 0	17.1 ± 0	18.52 ± 0	17.1 ± 0	18.52 ± 0	18.34 ± 0	17.1 ± 0	$\textbf{45.85} \pm \textbf{5.39}$	
30	17.1 ± 0	18.34 ± 0	17.1 ± 0	$\textbf{44.86} \pm \textbf{4.91}$					
40	17.1 ± 0	16.93 ± 0	17.1 ± 0	17.1 ± 0					
50	17.1 ± 0	16.93 ± 0	17.1 ± 0	17.1 ± 0					
60	17.1 ± 0	16.93 ± 0	17.1 ± 0	17.1 ± 0					
70	17.1 ± 0	16.93 ± 0	17.1 ± 0	17.1 ± 0					
80	17.1 ± 0	16.93 ± 0	17.1 ± 0	17.1 ± 0					
90	17.1 ± 0	16.93 ± 0	17.1 ± 0	17.1 ± 0					

Table 11 presents the classification accuracy of the MLP model, along with the baseline, benchmarks, and our mechanism, under SAFs, using the MNIST, FMNIST, and UCI-HAR datasets. The model with our mechanism classifies MNIST and FMNIST data samples more accurately than the model with the baseline and benchmarks, as shown in Section 6. We observe that the overall classification accuracy of the model decreases. This is because the precision for neurons' membrane potential and synaptic weights degrades due to low floating-point and bit width in the FPGA device, which damages the data stored in trained synaptic weights and neuronal activities.

B EFFICIENCY ANALYSIS BASED ON TIME/SPATIAL COMPLEXITY AND TIME/ENERGY CONSUMPTION

To demonstrate that our mechanism enhances the fault tolerance of SNN models without requiring complex algorithms, we measure the computational and spatial complexity of our mechanism and compare it to that of the benchmarks. Additionally, we measure the energy consumption of our mechanisms on the FPGA device since energy consumption is a significant advantage of SNNs that makes them suitable for neuromorphic device implementation.

B.1 COMPLEXITY ANALYSIS

We calculate the time and spatial complexities of the benchmarks and the proposed mechanism through experimental evidence that demonstrates the time and energy consumption of them in real devices.

B.1.1 TIME COMPLEXITY

We measure the time complexity of the benchmarks and the proposed mechanism. The following items present the time complexity of the benchmarks and the proposed mechanism.

- 1. **ECOC**: $T_{\text{ECOC}} = \Theta(EP^2n + NBCL)$.
- 2. **SoftSNN**: $T_{\text{SoftSNN}} = \Theta(NW)$.
 - 3. Routing: $T_{\text{route}} = \Theta \left(\sum_{\ell=1}^{L} C_{\text{out}}^{(\ell)} C_{\text{in}}^{(\ell)} k_{\ell}^{2} + \delta_{\text{swap}} \sum_{\ell=1}^{L} K_{\ell} \log K_{\ell} \right)$.
- 4. Astrocyte: $T_{Astro} = \Theta(P + NP)$.
- 5. FalVolt: $T_{\text{FalVolt}} = \Theta(M + N(M + fW))$.
- 6. LIFA: $T_{LIFA} = \Theta(P + NP)$
- 7. **Proposed**: $T_{\text{proposed, mlp}} = \Theta(N S [B + A + T + BTD]) + O(N_{\text{sal}} F_{\text{model}}^{\text{MLP}}), T_{\text{proposed, conv}} = \Theta(N S [B + A + T + BTD]).$

Notations of time complexity equations

- 1. **ECOC.** B: batch size; N: number of training/inference steps (batches processed); C: number of classes; E: number of extension code blocks; m: Hamming-code parameter; $n=2^m-1$: per code block length; $L=E\,n$: code length; P: candidate-pool size used in code book construction.
- 2. **Soft SNN.** $W = \sum_{i=1}^{M} P_i$: total number of trainable weights over the M layers; P_i : number of weights in layer i; $P_{\max} = \max_i P_i$: size of the largest layer.
- 3. **Routing.** L_{layers} : number of routed layers; $C_{\text{in}}^{(\ell)}$, $C_{\text{out}}^{(\ell)}$: input/output channels of layer ℓ ; k_{ℓ} : kernel size of layer ℓ (so $k_{\ell}^2 = 1$ for MLP/Linear); $K_{\ell} = \min\{C_{\text{in}}^{(\ell)}, C_{\text{out}}^{(\ell)}\}$: effective channel count for top-K matching; $\delta_{\text{swap}} \in \{0,1\}$: flag indicating whether sorting + channel-swap is enabled; W: total number of trainable weights across routed layers.
- 4. **Astrocyte.** N: number of batch iterations in an epoch; P: total number of trainable parameters over hooked layers; χ : output-channel chunk size used in the backward pass; $p_{\text{out}}^{\text{max}}$: maximum number of parameters associated with a single output channel (e.g., $9\,C_{\text{in}}$ for a 3×3 conv).
- 5. FalVolt. N: number of batch iterations in an epoch; W: total number of weights subject to potential fault mapping; M: number of spiking/protected modules whose thresholds or states are managed; $f \in [0,1]$: fraction of weights affected by faults (worst case f=1); $W_{\rm mask}$: number of stored fault-mask entries.

- 6. **LIFA.** N: number of batch iterations in an epoch; P: total number of trainable parameters across protected layers (for conv: $P = \sum_{\ell} C_{\mathrm{out}}^{(\ell)} C_{\mathrm{in}}^{(\ell)} k_{\ell}^2$; for linear: $k_{\ell}^2 = 1$); $C = \sum_{\ell} C_{\mathrm{out}}^{(\ell)}$: total output-channel count across protected layers.
- 7. **Proposed.** N: number of batch iterations in an epoch; $S = H \times W$: spatial size (pixels) per sample; B: batch size; D: input channels; T: number of time steps (fragments) per sample; A: number of orientation candidates; (optional only if the fault influence map (saliency and weight projection) is used) $N_{\rm sal}$: number of steps that compute saliency/backprop; $F_{\rm model}^{\rm MLP}$: per-step FLOPs of the MLP backbone under saliency.

With MNIST and FMNIST ($S = 28 \times 28 = 784$, D = 1), our fragmentation step scales as

$$T_{\text{proposed}} = \Theta(SBT) = \Theta(784BT), \tag{5}$$

while all weight-scanning benchmarks (Astrocyte and LIFA) based on astrocytes scale with the number of parameters:

$$T_{\text{scan}} = \Theta(W), \qquad T_{\text{FalVolt}} = \Theta(M + fW) \times \Theta(fW) \text{ (for non-vanishing } f),$$
 (6)

where W is the total trainable weights, M the number of spiking/protected modules, and $f \in [0, 1]$ the fraction of weights affected by faults. Hence, the decisive ratios are

$$\frac{T_{\text{proposed}}}{T_{\text{scan}}} \simeq \frac{784 \, B \, T}{W}, \qquad \frac{T_{\text{proposed}}}{T_{\text{FalVolt}}} \simeq \frac{784 \, B \, T}{f \, W}.$$
 (7)

For the CNNs in our setting (VGG-7/11/15, ResNet-18), W is in the multi-million range even on MNIST/FMNIST; with common batches/fragments ($B \in [64, 128]$, $T \in [2, 4]$) one has $W \gg 784~B~T$, so strictly $T_{\text{proposed}} < T_{\text{scanning}}$. The same conclusion holds against FalVolt for any fixed, non-negligible f (e.g., $f \ge 0.05$), since then $fW \gg 784~B~T$ in these networks, yielding $T_{\text{ours}} < T_{\text{FalVolt}}$ as well. ECOC differs in that its per-step cost is $T_{\text{ECOC}} = \Theta(B~C~L)$, giving

$$\frac{T_{\text{proposed}}}{T_{\text{ECOC}}} \approx \frac{STD}{CL} = \frac{784T}{CL},$$
 (8)

which is typically of the same order for $C{=}10$ and $L{\in}[64,256]$, while ECOC still incurs a one-off build of $\Theta(E\,P^2n)$. In summary, without the fault influence, our mechanism has strictly smaller per-step time complexity than all weight-scanning benchmarks under the MLP/VGG/ResNet models. It is competitive with (or smaller than) ECOC while avoiding the heavy one-time construction of ECOC. As shown in Appendix A.1 and A.2, our mechanism is not significantly dependent on the fault influence with CIFAR-10 and CIFAR-100 under VGG and ResNet models, indicating that our mechanism saves time by using only the complexity to make fragments under VGG and ResNet models with CIFAR-10 and CIFAR-100.

B.1.2 SPATIAL COMPLEXITY

We measure the spatial complexity of the benchmarks and the proposed mechanism.

- 1. ECOC: $S_{ECOC} = \Theta(P^2 + CL + BC)$.
- 2. **SoftSNN**: $S_{\text{SoftSNN}} = \Theta(W)$.
- 3. Routing: $S_{\text{route}} = \Theta(W)$.
- 4. Astrocyte: $S_{Astro} = \Theta(P + \chi p_{out}^{max})$.
- 5. FalVolt: $S_{\text{FalVolt}} = \Theta(M + W_{\text{mask}})$.
- 6. LIFA: $S_{LIFA} = \Theta(P+C)$.
- 7. **Proposed**: $S_{\text{proposed}} = \Theta(S [BTD + A])$.

Notations of spatial complexity equations

- 1. **ECOC.** C: number of classes; E: number of extension code blocks; m: Hamming parameter; $n=2^m-1$: code block length; $L=E\,n$: code length (size of the stored code book is $C\,L$); B: batch size (per-step logits buffer $B\,C$); P: candidate-pool size (build-time pairwise matrix P^2 gives peak memory).
- 2. **SoftSNN.** $W = \sum_{i=1}^{M} P_i$: total number of trainable weights over the M layers; $P_{\max} = \max_{i} P_i$: size of the largest layer (often dictates per-layer peak); $\mathcal{C} = \sum_{i=1}^{M} C_{\text{out},i}$: total number of output channels if per-channel thresholds are stored.
- 3. **Routing.** W: total number of trainable weights across routed layers (in-place operations keep the footprint parameter–scaled); $C_{\rm in}^{(\ell)}$: input channels of layer ℓ (small per-layer index/permutation buffers scale with $C_{\rm in}^{(\ell)}$ but are absorbed by W in big- Θ terms).
- 4. **Astrocyte.** P: total number of trainable parameters over hooked layers (CPU-side caches such as W_0 , inverse denominators, masks, q scale with P); χ : output-channel chunk size used on GPU during the backward pass; $p_{\text{out}}^{\text{max}}$: maximum number of parameters per output channel (e.g., $9 C_{\text{in}}$ for a 3×3 conv); the additional VRAM peak scales with $\chi p_{\text{out}}^{\text{max}}$.
- 5. FalVolt. W: total number of potentially fault-mapped weights (upper bound on parameter-scaled storage); M: number of protected/spiking modules (small bookkeeping state); $W_{\rm mask}$: number of stored fault-mask entries (persistent); typically $W_{\rm mask} \leq W$.
- 6. **LIFA.** P: total number of trainable parameters across protected layers (dominant persistent buffers: W_0 , inverse denominators, masks); $C = \sum_{\ell} C_{\text{out}}^{(\ell)}$: total number of output channels (per-channel EMA/state vectors).
- 7. **Proposed (fragmentation).** $S = H \times W$: per-sample spatial size (pixels); B: batch size; D: input channels; T: number of fragments per sample (fragment tensor B T D S dominates); A: number of orientation candidates (angle buffers A S).

With MNIST and FMNIST ($S=28\times 28=784,\,D=1$), the peak additional memory of the proposed mechanism scales as (Note that saliency map S is not essential for VGG and ResNet models.)

$$S_{\text{ours}} = \Theta(S[BTD + A]) = \Theta(784[BT + A]).$$
 (9)

where B is the batch size, T the number of fragments, and A the number of angle candidates. In contrast, weight–scan benchmarks that scan all parameters each step exhibit parameter–dominated footprints:

$$S_{\text{LIFA/Astro}} = \Theta(W)$$
 $S_{\text{SoftSNN}} = \Theta(P_{\text{max}}) \text{ (per-layer peak) } \lesssim \Theta(W).$ (10)

$$S_{\text{Routing}} = \Theta(W) \quad S_{\text{FalVolt}} = \Theta(M + W_{\text{mask}}) \le \Theta(W).$$
 (11)

where W is the total number of trainable weights, $P_{\rm max}$ the size of the largest layer's weight tensor, M the number of protected/spiking modules, and $W_{\rm mask}$ the number of stored fault–mask entries. Therefore, for typical MNIST/FMNIST settings (e.g., $B \in [64, 128], T \in [2, 4], A \leq 180$) and CNN backbones (VGG-7/11/15, ResNet-18) with W in the multi–million range,

$$\frac{S_{\rm proposed}}{S_{\rm LIFA/Astro}} \approx \frac{784 \left[B\,T + A\right]}{W} \, \ll \, 1, \qquad \frac{S_{\rm proposed}}{S_{\rm FalVolt}} \, \approx \, \frac{784 \left[B\,T + A\right]}{M + W_{\rm mask}} \, \ll \, 1. \tag{12}$$

and similarly $S_{\text{ours}} \ll S_{\text{Routing}}$ and $S_{\text{ours}} \lesssim S_{\text{SoftSNN}}$ whenever P_{max} is large (as in VGG/ResNet). ECOC is different: per–step it stores only the code book and logits,

$$S_{\text{ECOC, step}} = \Theta(CL + BC) \quad (C=10 \text{ on MNIST/FMNIST}, L = \mathcal{O}(10^2)).$$
 (13)

which is often smaller than S_{proposed} on these datasets; however, ECOC incurs a one–time build peak of $\Theta(P^2)$ (candidate–pair matrix) that can dominate transient memory. In summation, without the fault influence map, our method achieves strictly smaller parameter space complexity than all weight–scan benchmarks (LIFA, Astrocyte, Soft SNN, Routing, FalVolt), while remaining competitive with ECOC apart from its negligible per–step footprint but heavy one–off construction of ECOC.

B.2 TRAINING TIME

 We measure the training time of the baseline, benchmarks, and the proposed mechanism using the MLP (MNIST), VGG-7 (CIFAR-10), ResNet-18 (CIFAR-100), and ResNet-34 (Tiny-ImageNet) models with 2 time steps. We train the models on a workstation with an Nvidia GeForce RTX 4080 GPU with Ubuntu 24.04.

Table 12: Various models' training time (sec) in a 95% confidence interval with the baseline, benchmarks, and proposed mechanism on a workstation under SAFs with a fault ratio of 0.5 and 2 time steps.

	Baseline	ECOC	SoftSNN	Routing	Astrocyte	FalVolt	LIFA	Proposed
	Dascillic	LCOC	301131414	Kouting	Astrocyte	1 ai voit	LIIA	1 Toposcu
MLP	$\textbf{193.84} \pm \textbf{2.71}$	197.62 ± 2.85	196.51 ± 2.56	198.29 ± 3.05	288.75 ± 4.51	201.47 ± 3.21	293.86 ± 4.77	205.24 ± 4.23
VGG-7	291.16 ± 3.57	296.91 ± 3.8	294.34 ± 3.65	298.81 ± 4.01	351.82 ± 5.27	303.53 ± 4.26	356.74 ± 5.53	310.38 ± 5.18
ResNet-18	382.53 ± 4.12	385.61 ± 4.03	384.77 ± 4.53	396.54 ± 4.68	721.97 ± 6.28	408.9 ± 4.94	724.62 ± 6.09	413.32 ± 4.5
ResNet-34	4259.57 ± 18.62	4304.4 ± 20.11	4298.46 ± 21.39	4350.83 ± 25.75	8005.37 ± 36.21	4317.89 ± 26.04	8154.17 ± 32.83	4392.13 ± 26.51

Our mechanism consumes significantly less training time than weight-scanning approaches based on astrocytes (Astrocyte and LIFA), as we demonstrate that our mechanism definitely consumes less time than the astrocyte-based approaches due to their less complexity. Unlike these approaches, the training time of the models with our mechanism does not increase significantly as the complexity of the models and datasets increases. The model with ours also consumes comparable training time to that of ECOC, SoftSNN, Routing, and FalVolt. This evaluation result shows that our mechanism does not severely inflate the burden on training time.

B.3 ENERGY CONSUMPTION ON THE REAL FPGA DEVICE

We measure the energy consumption of the model with the baseline, benchmarks, and proposed mechanism on the FPGA device during testing. Table 13 exhibits the energy consumption of the model with MNIST and FMNIST on the FPGA device using 2 time steps.

Table 13: The MLP models' energy consumption (mJ) in a 95% confidence interval with the baseline, benchmarks, and proposed mechanism on the real FPGA hardware with two time steps.

	Baseline	ECOC	SoftSNN	Routing	Astrocyte	FalVolt	LIFA	Proposed
MNIST	85.31 ± 1.05	88.76 ± 1.27	86.23 ± 1.18	90.44 ± 1.31	150.72 ± 1.61	95.15 ± 1.23	165.69 ± 1.59	$\textbf{67.16} \pm \textbf{0.82}$
FMNIST	87.19 ± 1.36	90.54 ± 1.57	88.68 ± 1.43	92.82 ± 1.51	156.08 ± 1.99	98.23 ± 1.72	168.33 ± 1.93	$\textbf{78.37} \pm \textbf{0.95}$

The MLP model with our mechanism exhibits the least energy consumption among the MLP models on the real FPGA device. This is because our mechanism shrinks the size of the data samples through fragmentation, and the probability of spike occurrence declines since the number of non-zero pixels decreases during fragmentation, as mentioned in Subsection 5.3. This effect enables the model with our mechanism to consume less energy than the models with all benchmarks, despite our mechanism having higher time complexity and consumption than some benchmarks. However, the benchmarks increase the complexity of the decoding (ECOC), keep neurons' activation frequent (SoftSNN and Routing), utilize the astrocyte module to activate non-faulty synapses (Astrocyte and LIFA), and incorporate additional learnable parameters to adjust neuronal activities (Falvolt). Therefore, the MLP models with the benchmark require more energy than ours.

C DETAILED MATHEMATICAL EXPLANATION OF THE MOTIVATION STUDY

We demonstrate how synaptic faults ruin the usable learning capacity of SNN models mathematically.

C.1 SETUP AND NOTATION

 Consider a spiking neuron with membrane potential $V_t \in \mathbb{R}$, threshold $\vartheta \in \mathbb{R}$, and spike output $K_t \in \{0,1\}$. During training, we replace the Heaviside step H by a surrogate $\sigma : \mathbb{R} \to [0,1]$ so that $K_t \approx \sigma(V_t - \vartheta)$ and σ' is used in backpropagation. Let the *surrogate gradient corridor width* be $\delta > 0$ such that $\sigma'(u) \approx 0$ whenever $|u| > \delta$. For a feedforward pre-activation at layer ℓ and time t,

$$z_t^{(\ell)} = W^{(\ell)} K_t^{(\ell-1)} + b^{(\ell)}$$
 (vector form), (14)

and for a single neuron with input $x \in \mathbb{R}^d$ and weights $w \in \mathbb{R}^d$ we write $z = w^\top x + b$. For an LIF neuron, we use

$$V_t = \alpha V_{t-1} + z_t - \vartheta K_{t-1}, \qquad K_t \approx \sigma(V_t - \vartheta), \qquad \alpha \in (0, 1).$$
 (15)

C.2 SURROGATE GRADIENT CORRIDOR

Let $u := z - \vartheta$. Many arctangent surrogates used in SNNs have a backward derivative of the rational form

$$\phi'(u) = \frac{A}{1 + (\beta u)^2}, \qquad A > 0, \ \beta > 0, \tag{16}$$

which yields, for a target gradient floor $\gamma \in (0, A)$, the corridor

$$C_{\gamma} := \{ u : \phi'(u) \ge \gamma \} = [-\delta(\gamma), \delta(\gamma)], \quad \delta(\gamma) = \frac{1}{\beta} \sqrt{\frac{A}{\gamma} - 1}.$$
 (17)

Additionally, we derive the surrogate gradient corridor of the arctangent function, which is widely used as a surrogate gradient function for LIF neurons.

Let $u := z - \vartheta$ and consider the arctangent surrogate derivative

$$\phi'(u) = \frac{\alpha_s}{\pi (1 + (\alpha_s u)^2)}, \qquad \alpha_s > 0.$$
(18)

For a target gradient floor $\gamma \in (0, \alpha_s/\pi)$, define the corridor(Li et al., 2024a; Zenke & Vogels, 2021; Shrestha & Orchard, 2018)

$$C_{\gamma} := \{ u : \phi'(u) \ge \gamma \} = [-\delta, \delta], \qquad \delta(\gamma) = \frac{1}{\alpha_s} \sqrt{\frac{\alpha_s}{\pi \gamma} - 1}.$$
 (19)

whose peak is $A = 1/\pi$. Setting the gradient floor to $r = fA = f/\pi$ yields the corridor half-width

$$\delta(f) = \frac{2}{\pi\alpha} \sqrt{\frac{1}{f} - 1}.$$
 (20)

In practice, we initialize $\alpha=2$ and f=0.2 (thus $r\approx0.0637$), then adapt f per layer using mini-batch membrane statistics so that the corridor covers a target score p of the observed $U-\vartheta$ distribution: with $\hat{\sigma}_{\ell}$ the running standard deviation and z_p the normal quantile for score p, we set $\delta(f_{\ell})\approx z_p\hat{\sigma}_{\ell}$, i.e (Zenke & Vogels, 2021; Wang et al., 2023; Che et al., 2022; Lian et al., 2023).

$$f_{\ell} = \frac{1}{1 + \left(\frac{\pi\alpha}{2} z_{p} \hat{\sigma}_{\ell}\right)^{2}}.$$
 (21)

This keeps most samples within the high-gradient band while avoiding an overly narrow corridor.

Remark 1 (Mapping to implementation). For the common parameterization $\phi'(u) = \frac{\alpha_s}{\pi(1+(\alpha_s u)^2)}$, one has $A = \alpha_s/\pi$ and $\beta = \alpha_s$. For the SpikingJelly ATan $\phi'(u) = \frac{\alpha/2}{1+(\frac{\pi\alpha}{2}u)^2}$, one has $A = \alpha/2$ and $\beta = \pi\alpha/2$ (Fang et al., 2023). Both are instances of equation 16, so equation 17 applies verbatim.

C.3 FAULT MODELING

We consider synaptic faults that perturb parameters and/or inputs:

$$w \mapsto w + \Delta w, \qquad x \mapsto x + \Delta x,$$
 (22)

where $\Delta w, \Delta x$ may be sparse (e.g., stuck-at-0/1 (SA0/SA1) at a subset of synapses) or dense (e.g., analog drift). The post-fault pre-activation is

$$z' = (w + \Delta w)^{\top} (x + \Delta x) + b = z + \Delta z, \quad \Delta z = \underbrace{\Delta w^{\top} x}_{\text{param fault}} + \underbrace{w^{\top} \Delta x}_{\text{input fault}} + \underbrace{\Delta w^{\top} \Delta x}_{\text{higher-order}}.$$
 (23)

By Cauchy-Schwarz,

$$|\Delta z| \le ||x||_2 ||\Delta w||_2 + ||w||_2 ||\Delta x||_2 + ||\Delta w||_2 ||\Delta x||_2. \tag{24}$$

SA0 on an input line j is modeled by $(\Delta x)_j = -x_j$; SA1 by $(\Delta x)_j = c - x_j$ for a fixed logic level c. Bit/weight stuck faults are included in Δw .

C.4 FROM FAULTS TO SATURATION

At time t, the only instantaneous change from a synaptic fault is $z_t \mapsto z_t + \Delta z_t$, hence

$$V_t' = \alpha V_{t-1} + (z_t + \Delta z_t) - \vartheta K_{t-1} = V_t + \Delta z_t, \quad \Rightarrow \quad V_t' - \vartheta = (V_t - \vartheta) + \Delta z_t. \quad (25)$$

Lemma 1 (Corridor escape: sufficient conditions). Let $a_t := V_t - \vartheta$ and suppose $|a_t| \le \delta$ (pre-fault state inside the corridor).

- 1. (Sign-aligned escape) If $a_t \Delta z_t \geq 0$ and $|\Delta z_t| \geq \delta |a_t|$, then $|a_t + \Delta z_t| \geq \delta$, hence $\sigma'(V_t' \vartheta) \approx 0$ at time t.
- 2. (Sign-agnostic escape) Regardless of the sign of Δz_t , if $|\Delta z_t| > \delta + |a_t|$, then $|a_t + \Delta z_t| > \delta$.

Proof. (1) If $a_t \Delta z_t \ge 0$ then $|a_t + \Delta z_t| = ||a_t| + |\Delta z_t|| \ge \delta$ when $|\Delta z_t| \ge \delta - |a_t|$. (2) By the reverse triangle inequality, $|a_t + \Delta z_t| \ge ||\Delta z_t| - |a_t|| > \delta$.

C.5 EXPECTED GRADIENT BOUND FOR A SINGLE NEURON

Let $g_t := \partial \mathcal{L}/\partial S_t$ and suppose $0 \le \sigma'(u) \le C_\sigma \mathbf{1}\{|u| \le \delta\}$. Then,

$$\left\| \frac{\partial \mathcal{L}}{\partial w} \right\| = \left\| \sum_{t=1}^{T} g_t \, \sigma'(V_t' - \vartheta) \, x_t \right\| \leq C_\sigma \sum_{t=1}^{T} \|g_t\| \, \|x_t\| \, \mathbf{1}\{|V_t' - \vartheta| \leq \delta\}. \tag{26}$$

Taking expectations and using Cauchy-Schwarz yields the model-free bound

$$\mathbb{E}\left\|\frac{\partial \mathcal{L}}{\partial w}\right\| \leq C_{\sigma} \sum_{t=1}^{T} \left(\mathbb{E}\left[\|g_{t}\|^{2} \|x_{t}\|^{2}\right]\right)^{1/2} \cdot \mathbb{P}(|V'_{t} - \vartheta| \leq \delta)^{1/2}. \tag{27}$$

Under a mild independence/mixing assumption between $||g_t|| ||x_t||$ and the corridor event, one may write the simpler scaling

$$\mathbb{E}\left\|\frac{\partial \mathcal{L}}{\partial w}\right\| \lesssim C_{\sigma} \sum_{t=1}^{T} \mathbb{E}\left[\|g_{t}\| \|x_{t}\|\right] p_{t}, \qquad p_{t} \coloneqq \mathbb{P}(|V'_{t} - \vartheta| \le \delta).$$
 (28)

C.6 DEPTH- AND TIME-WISE COMPOUNDING

For a parameter in layer ℓ , a generic backpropagation path contains factors $\sigma'(V_t^{(j)} - \vartheta)$ for $j \leq \ell$ and relevant t. Bounding σ' by indicators,

$$|\Pi| \le C_{\sigma}^{N_{\Pi}} \prod_{j,t: \gamma_{j,t}=1} \mathbf{1}\{|V_{t}^{(j)} - \vartheta| \le \delta\}, \qquad N_{\Pi} = \sum_{j,t} \gamma_{j,t}.$$
 (29)

Taking expectations gives

$$\mathbb{E}\left|\Pi\right| \leq C_{\sigma}^{N_{\Pi}} \mathbb{P}\left(\bigcap_{j,t:\,\gamma_{j,t}=1} \{|V_{t}^{(j)} - \vartheta| \leq \delta\}\right). \tag{30}$$

A conservative bound is

$$\mathbb{E}\left|\Pi\right| \leq C_{\sigma}^{N_{\Pi}} \min_{j,t: \gamma_{j,t}=1} p_{j,t}, \qquad p_{j,t} \coloneqq \mathbb{P}\left(\left|V_{t}^{(j)} - \vartheta\right| \leq \delta\right). \tag{31}$$

If corridor events are approximately independent (or satisfy a weak-mixing condition), then

$$\mathbb{E}\left|\Pi\right| \leq C_{\sigma}^{N_{\Pi}} \prod_{j,t: \gamma_{j,t}=1} p_{j,t} \leq C_{\sigma}^{N_{\Pi}} (p^{\star})^{N_{\Pi}}, \qquad p^{\star} \coloneqq \sup_{j,t} p_{j,t}, \tag{32}$$

exhibiting exponential attenuation as N_{Π} grows.

C.7 FIRST LAYER SENSITIVITY IN MLP

For the first layer (vector form) with $z^{(1)} = W^{(1)}x + b^{(1)}$ and perturbations $(\Delta W^{(1)}, \Delta x)$,

$$\|\Delta z^{(1)}\| \le \|\Delta W^{(1)}\|_{\text{op}} \|x\|_2 + \|W^{(1)}\|_{\text{op}} \|\Delta x\|_2 + \|\Delta W^{(1)}\|_{\text{op}} \|\Delta x\|_2, \tag{33}$$

so sizeable input/weight faults directly shift $z^{(1)}$ without any preceding contraction, shrinking corridor occupancy in deeper layers via equation 30–equation 32.

C.8 SUFFICIENT CONDITION FOR GRADIENT COLLAPSE

Define $p_{j,t}$ as above and let \mathcal{G} be the multiset of "corridor gates" along dominant backpropagation paths with size N_* . If a fraction $\rho \in (0,1]$ of gates satisfy $p_{j,t} \leq \varepsilon \ll 1$, then

$$\mathbb{E} |\Pi| \leq \begin{cases} C_{\sigma}^{N_*} \varepsilon^{\rho N_*}, & \text{under independence/mixing,} \\ C_{\sigma}^{N_*} \varepsilon, & \text{(conservative, no independence).} \end{cases}$$
(34)

Either case shows attenuation; the independent/mixing case yields exponential decay in depth \times time.

C.9 EFFECTIVE BIAS INTERPRETATION FOR SA0/SA1 of SAFS

For SA1 on a subset \mathcal{J} of input lines with logic level c,

$$\Delta z = w^{\top} \Delta x = \sum_{j \in \mathcal{J}} w_j (c - x_j) = c \sum_{j \in \mathcal{J}} w_j - \sum_{j \in \mathcal{J}} w_j x_j, \tag{35}$$

acting as an additive bias shift plus removal of signal terms. Persistent shifts displace V_t away from ϑ across time steps, driving down corridor occupancy $p_{j,t}$ and compounding the bottleneck via equation 32.

C.10 SUMMARY

Synaptic faults induce a pre-activation shift Δz decomposed in equation 23 and bounded in equation 24. When $|\Delta z|$ is large relative to the corridor width δ , Lemma 1 ensures $|V_t - \vartheta| > \delta$ so $\sigma'(V_t - \vartheta) \approx 0$. The expected gradient is then attenuated proportionally to (at least) $\sqrt{p_{j,t}}$ per time step equation 27; under independence/mixing, it scales with $p_{j,t}$ equation 28. Across layers and time steps, this attenuation multiplies equation 32, producing the *bottleneck problem*, with the first layer of MLP especially vulnerable by equation 33.

D NEAR-OPTIMALITY OF THE PROPOSED MECHANISM

We show why our solution is the near-optimal solution to improve the fault tolerance of SNNs in this section.

D.1 SETUP AND NOTATION

We consider inputs $x \in \mathbb{R}^n$ and a fixed number of stripes (1D profiles in Section 5) $T \in \mathbb{N}$. Indices are $i \in \{1,\ldots,n\}$ and stripes are $t \in \{1,\ldots,T\}$. A stripe partition is represented by binary masks $M_t(i) \in \{0,1\}$ that satisfy $\sum_{t=1}^T M_t(i) = 1$ for every i, and contiguity is taken with respect to a one-dimensional scan order of the indices induced by an angle θ in a finite set $\Theta \subset [0,\pi)$. Given any nonnegative vector $s \in \mathbb{R}^n_+$, the load of stripe t is the linear functional $S_t(s) = \sum_i s_i M_t(i)$; we also write the total mass $U(s) = \sum_i s_i$, the per-stripe mean $\mu(s) = U(s)/T$, and the element-wise maximum $m(s) = \max_i s_i$. For vectors $a, b \in \mathbb{R}^n$, the inner product is $\langle a, b \rangle = \sum_i a_i b_i$ and $\|v\|_p$ denotes the ℓ_p norm; the Hadamard product is $a \odot b$.

$$S_t(s) := \sum_{i=1}^n s_i M_t(i), \qquad U(s) := \sum_{i=1}^n s_i, \qquad \mu(s) := \frac{U(s)}{T}, \qquad m(s) := \max_i s_i.$$
 (36)

Given trained weights $w \in \mathbb{R}^n$ and a fault/perturbation Δw , we set $\widehat{w} := w + \Delta w$ and restrict the input to stripe t by $x_t := x \odot M_t$. The (stripe) pre-activation is

$$z_t := \langle \widehat{w}, x_t \rangle, \qquad x_t := x \odot M_t, \qquad \widehat{w} := w + \Delta w.$$
 (37)

We denote by $z^* > 0$ the corridor threshold, i.e., the largest value for which the chosen surrogate derivative $\phi'(z)$ remains in its effective (non-vanishing) regime for all $|z| \le z^*$. To construct stripes, we employ an implementable importance map $I \in \mathbb{R}^n_+$ and assume a two-sided calibration with respect to the ideal per-index load $u_i := |u_i| |x_i|$: there exist constants $0 < c_- \le 1 \le c_+$ such that

$$c_{-}|w_{i}||x_{i}| \leq I_{i} \leq c_{+}|w_{i}||x_{i}|, \qquad i = 1, \dots, n.$$
 (38)

When I = u one has $c_- = c_+ = 1$. The quantile (greedy) stripes used in the paper are obtained by scanning indices in the chosen order and inserting a cut whenever the cumulative load with respect to I first exceeds integer multiples of $\mu(I)$, producing T contiguous fragments.

D.2 FAULT MODELS AND A BASIC UPPER BOUND

We consider the following three fault models, mentioned in Section 2.

SAFs: Some synapses are permanently stuck at G_{\min} or G_{\max} so the implemented weight becomes w_i' (e.g., SA0/SA1). Let $\Delta w_i := w_i' - w_i$ and assume $\|\Delta w\|_{\infty} \leq \varepsilon_{\text{SAF}}$. Then for any stripe t (Boyd & Vandenberghe, 2004),

$$|z_{t}| = |\langle w + \Delta w, x_{t} \rangle| \leq \underbrace{\sum_{i} |w_{i}| |x_{i}| M_{t}(i)}_{S_{t}(u)} + \varepsilon_{\text{SAF}} \underbrace{\sum_{i} |x_{i}| M_{t}(i)}_{S_{t}(|x|)} = S_{t}(u) + \varepsilon_{\text{SAF}} S_{t}(|x|).$$
(39)

 RWFs: Each coordinate experiences an independent, mean-zero, bounded (or sub-Gaussian) perturbation Δw_i . If $|\Delta w_i| \leq b$ and the Δw_i are independent, then for any $\tau > 0$ and stripe t,

$$\Pr(|\langle \Delta w, x_t \rangle| > \tau) \le 2 \exp\left(-\frac{\tau^2}{2b^2 \|x_t\|_2^2}\right), \tag{40}$$

so equalizing $S_t(|x|) = ||x_t||_1$ across t uniformly tightens the tail bound (sub-Gaussian and Hoeffding) (Hoeffding, 1963).

CEFs: Wiring errors apply a linear transformation to the input so that $z_t = w^{\top}(Ax_t)$. This is equivalent to using the effective weight $w' := A^{\top}w$, i.e., $\Delta w^{(c)} := (A^{\top} - I)w$. If $\|\Delta w^{(c)}\|_{\infty} \le \varepsilon_{\text{CEF}}$, then

$$|z_t| = |\langle w + \Delta w^{(c)}, x_t \rangle| \le S_t(u) + \varepsilon_{\text{CEF}} S_t(|x|).$$
(41)

A permutation fault A = P is a special case; taking $\varepsilon_{\text{CEF}} = \|(P^\top - I)w\|_{\infty}$ yields the same bound (Boyd & Vandenberghe, 2004).

D.3 CALIBRATION: ALIGNING THE IMPORTANCE MAP WITH THE EFFECTIVE PER-INDEX LOAD

We formalize the requirement that the implementable importance I should approximate u within stripe-wise sums.

Assumption 1 (Two-sided calibration.). There exist constants $0 < c_{-} \le 1 \le c_{+}$ such that for all indices i,

$$c_{-}u_{i} \leq I_{i} \leq c_{+}u_{i}. \tag{42}$$

Lemma 2 (Calibration). Under the assumption (Two-sided calibration), for any stripe partition,

$$S_t(u) \le \frac{1}{c_-} S_t(I), \qquad \mu(u) \le \frac{1}{c_-} \mu(I), \qquad m(u) \le \frac{1}{c_-} m(I).$$
 (43)

Proof. From $u_i \leq I_i/c_-$, sum over i in stripe t. Similar for totals and maxima.

D.4 QUANTILE STRIPES ARE ADDITIVELY NEAR-OPTIMAL (CONTIGUOUS CASE)

Fix a nonnegative sequence a_1, \ldots, a_n obtained by scanning the image along any 1D order (e.g., the θ -scan used in the main text). Let $U(a) = \sum_i a_i$ and target mean $\mu(a) = U(a)/T$. Define the quantile (greedy) contiguous partition in Subsections 5.2 and 5.3 by sweeping from left to right and cutting whenever the cumulative sum first exceeds multiples of $\mu(a)$, producing T contiguous stripes.

Lemma 3 (Additive bound for greedy quantiles). Let $m(a) := \max_i a_i$. Then the greedy quantile partition satisfies

$$\max_{t \le T} S_t(a) \le \mu(a) + m(a). \tag{44}$$

Moreover, any contiguous partition must have $\max_t S_t(a) \ge \mu(a)$; hence, the greedy partition is a +m(a)-additive approximation to the optimal contiguous partition.

Proof. Each of the first T-1 stripes stops at the first index that causes the running sum to exceed $\mu(a)$. The overshoot over $\mu(a)$ is therefore at most the last included element, i.e., $\leq m(a)$. Hence every one of the first T-1 stripes has load in $(\mu(a), \, \mu(a) + m(a)]$. The final stripe has the remaining mass $U(a) - \sum_{t=1}^{T-1} S_t(a) \leq \mu(a)$. Thus, the maximum stripe load is at most $\mu(a) + m(a)$. The lower bound $\geq \mu(a)$ holds by a pigeonhole argument.

Theorem 1 (Near–optimality for u via quantiles on I). Construct stripes by greedy quantiles on the calibrated importance I. Under Assumption 1:

$$\max_{t} S_t(u) \leq \frac{1}{c_-} \left(\mu(I) + m(I) \right). \tag{45}$$

If I=u (so $c_-=c_+=1$), the greedy partition achieves $\max_t S_t(u) \leq \mu(u) + m(u)$, i.e., a+m(u) additive approximation to the optimal contiguous value. For a calibrated I with Assumption I, we have $\max_t S_t(u) \leq \frac{1}{c_-} \left(\mu(I) + m(I) \right)$. Translating this bound to the u-optimum introduces a calibration-dependent drift via $\mu(I) \in [c_-\mu(u), c_+\mu(u)]$, so the additive gap to the optimal contiguous value is at most $\frac{1}{c_-} m(I) + \frac{c_-}{c_-} - 1 \mu(u)$. The baseline $\mu(u)$ is at most $\frac{1}{c_-} \mu(I)$. Proof. Apply Lemma 3 with a=I to obtain $\max_t S_t(I) \leq \mu(I) + m(I)$. Then use Lemma 2: $S_t(u) \leq S_t(I)/c_-$.

Remark 2 (Direct partitioning condition). (i) If one directly partitions using a = u, then $c_- = 1$ and the bound gives $\max_t S_t(u) \le \mu(u) + m(u)$. (ii) The proof does not assume the particular 1D order beyond contiguity; the order may be induced by any scan (e.g., the θ -parameterization used to define stripes).

D.5 CORRIDOR PRESERVATION: SUFFICIENT CONDITIONS

Under SAFs and CEFs, combining the bounds with Theorem 1 yields a closed-form uniform bound on $|z_t|$:

$$|z_t| \le \underbrace{\frac{1}{c_-} \left(\mu(I) + m(I)\right)}_{\text{from } u} + \underbrace{\varepsilon S_t(|x|)}_{\text{fault term}}.$$
 (46)

The stripes are constructed by greedy quantiles on I (not on |x|). Let $X_{\max} := \max_{t \leq T} \operatorname{St}(|x|)$, computed on the *same* I-quantile partition. A simple sufficient condition for staying within the corridor is then

$$\frac{1}{c_{-}} \big\{ \mu(I) + m(I) \big\} + \varepsilon X_{\max} \leq z^{\star}, \qquad X_{\max} := \max_{t \in [T]} S_{t}(|x|) \text{ (computed on the } same \text{ I-quantile partition)}. \tag{47}$$

Remark 3 (Optional (co-monotone scan)). *If along the 1D scan used to build the I-quantile partition the sequences I and* |x| *are approximately co-monotone—so that applying Lemma 3 to* |x| *is justified—then*

$$X_{\text{max}} \le \mu(|x|) + m(|x|).$$
 (48)

In this case, a convenient sufficient condition is

$$\frac{1}{c} \left\{ \mu(I) + m(I) \right\} + \varepsilon \left(\mu(|x|) + m(|x|) \right) \le z^{\star}. \tag{49}$$

Under RWFs, Hoeffding's tail (and a union bound) implies that with probability at least $1 - 2T \exp\{-\tau^2/(2b^2 \max_t \|x_t\|_2^2)\}$, all stripes satisfy $|\langle \Delta w, x_t \rangle| \le \tau$ (Hoeffding, 1963). Thus, the (random) bound analogous to the above holds with the deterministic term $\varepsilon S_t(|x|)$ replaced by τ , chosen at the desired confidence level.

D.6 ON THE GINI OBJECTIVE (PRIMARY SURROGATE FOR MIN-MAX LOAD)

We treat minimizing the Gini coefficient of the 1D projection of I as a **primary surrogate** for suppressing the worst-case stripe load. Recall that G(S) equals one-half of the relative mean absolute difference and is equivalent to the Lorenz-based definition. Hence, it directly reduces pairwise dispersion. The next proposition turns this dispersion control into a deviation bound that is *linear* in G(S) (Yitzhaki & Schechtman, 2013).

Proposition 1 (Gini \Rightarrow deviation bound). Let $S \in \mathbb{R}^T_+$ with mean μ and Gini coefficient G(S). Then

$$\max_{t} |S_{t} - \mu| \leq \frac{1}{T} \sum_{i,j} |S_{i} - S_{j}| = 2T \,\mu \,G(S).$$

1890 Proof. By $\sum_{j}(S_t-S_j)=T(S_t-\mu)$ and the triangle inequality, $|S_t-\mu|=\frac{1}{T}\big|\sum_{j}(S_t-S_j)\big|\leq \frac{1}{T}\sum_{j}|S_t-S_j|$. Summing over t and taking the maximum yields $\max_{t}|S_t-\mu|\leq \frac{1}{T}\sum_{i,j}|S_i-S_j|$. Since $\sum_{i,j}|S_i-S_j|=2T^2\mu G(S)$, the claim follows.

Combine Proposition 1 with the additive near-optimality bound for contiguous quantile stripes (Lemma/Theorem: $\max_t S_t \leq \mu + m$ for the greedy split). Minimizing G(S) tightens $\max_t |S_t - \mu|$ and thus reduces $\max_t S_t$ under the same partition, making the corridor constraint strictly easier to satisfy. In short: Gini $\downarrow \Rightarrow$ pairwise dispersion $\downarrow \Rightarrow$ deviation $\downarrow \Rightarrow$ min-max load \downarrow .

D.7 COMPUTING THE SCAN/STRIPES

Let $\Theta \subset [0,\pi)$ denote a finite set of scan angles (or any family of 1D orders). For a fixed order, the greedy quantile partition runs in linear time. If one wishes to *search* over Θ , evaluate the objective $\max_t S_t(I)$ for each order and pick the best; since the objective only changes at permutation "event points", coarse uniform sampling of Θ is typically sufficient in practice. When an exact optimum over contiguous partitions is desired for a fixed order, classical Dynamic Programming (DP) or feasibility—check with binary search finds $\min_{\text{contig}} \max_t S_t(I)$ in polynomial time; our greedy rule is the simple additive-approximate alternative used in the paper (Skiena, 2008).

D.8 SUMMARY

For any calibrated I, greedy quantile stripes achieve the near-optimality bound above; if I=u, the achieved maximum load is within +m(u) of the contiguous optimum for u. Under the SAFs and CEFs, the closed-form sufficient condition ensures $|z_t| \leq z^*$ for all stripes, preventing gradient collapse; under RWFs, the analogous high-probability statement follows from the sub-Gaussian tail bound. The constants involved are the calibration c_- , the fault radius ε (or (b,τ) in the probabilistic model), and the observable statistics $\mu(\cdot)$ and $m(\cdot)$.

E FAULT-TOLERANCE CAPACITY PREDICTION OF OUR MECHANISM

In this section, we analyze the fault-tolerance capacity of our mechanism under the SAF, RWF, and CEF injection using arctangent as a surrogate gradient function. The boundary of synaptic weights is [-1,1] (Le Gallo et al., 2023; Lammie et al., 2022).

E.1 SETUP AND NOTATION

Let $w \in [-1,1]^N$ be the clean weight vector, $K := \|w\|_2$, and let a fraction $\rho \in [0,1]$ of synapses be faulty. All results below are per layer and can be applied layer-wise. Please refer to the derivation of the surrogate gradient (arctangent) corridor in Subsection C.2 while reading our paragraphs on capacity calculation.

With dynamic fragmentation and per-fragment RMS normalization to $\|\tilde{x}_t\|_2 = \alpha_n$,

$$|u_t| = |\hat{w}^\top \tilde{x}_t + b - \vartheta| \le ||\hat{w}||_2 \alpha_n + m, \quad m := |b - \vartheta|.$$
 (50)

Hence it suffices that $\|\hat{w}\|_2 \leq B$, where

$$B := \frac{\delta(\gamma) - m}{\alpha_n} \quad \text{(requires } \delta(\gamma) > m\text{)}. \tag{51}$$

E.2 CAPACITY UNDER SAFS

Under SA0, we replace faulty entries by -1; under SA1 by +1. In either case $|\hat{w}_i| = 1$ on faulty indices, so

$$\|\hat{w}\|_{2}^{2} = \|w\|_{2}^{2} - \sum_{i \in F} w_{i}^{2} + \rho N \cdot 1 \le K^{2} + \rho N,$$
(52)

where the inequality is the *deterministic worst-case* (we drop the nonnegative subtraction term). Therefore, a sufficient condition to remain inside the corridor is

$$K^2 + \rho N \leq B^2 \implies \rho^{\star}_{\mathrm{SA}\pm, \ worst} = \frac{B^2 - K^2}{N}$$
 (clipped to $[0, 1]$).

(53)

(54)

(55)

If faulty indices are drawn uniformly at random (independent of
$$w$$
), then $\mathbb{E}\big[\sum_{i\in F} w_i^2\big] = \rho K^2$ and

$$\mathbb{E}\|\hat{w}\|_{2}^{2} = (1-\rho)K^{2} + \rho N \cdot 1 = K^{2} + \rho (N - K^{2}),$$

E.3 CAPACITY UNDER RWFS

On faulty indices, $\hat{w}_i = w_i + \varepsilon_i$ with $\mathbb{E}[\varepsilon_i] = 0$ and $\text{Var}(\varepsilon_i) = \sigma_w^2$. Independence yields $\mathbb{E}\|\hat{w}\|_2^2 = 0$ $K^2 + \rho N \sigma_w^2$, so

 $\rho_{\text{SA}\pm, \ exp}^{\star} = \frac{B^2 - K^2}{N - K^2}$ (clipped to [0, 1]).

$$\rho_{\text{RWF, }exp}^{\star} = \frac{B^2 - K^2}{N \sigma_w^2} \text{ (clipped to } [0, 1]).$$
 (56)

A high-probability version follows from sub-Gaussian concentration by replacing $N\sigma_w^2$ with an upper-tail bound.

E.4 CAPACITY UNDER CEFS

A fraction ρ of entries are replaced by i.i.d. U[a,b] and then frozen. Let $\mu_f = \frac{a+b}{2}$ and $\sigma_f^2 = \frac{(b-a)^2}{12}$ so that $\mathbb{E}[\hat{w}_i^2] = \mu_f^2 + \sigma_f^2$. If faulty indices are random (independent of w),

$$\mathbb{E}\|\hat{w}\|_{2}^{2} = (1-\rho)K^{2} + \rho N(\mu_{f}^{2} + \sigma_{f}^{2}) \Rightarrow \rho_{\text{CEF, }exp}^{\star} = \frac{B^{2} - K^{2}}{N(\mu_{f}^{2} + \sigma_{f}^{2}) - K^{2}} \text{ (clipped to } [0, 1]).$$
(57)

For the common symmetric case U[-1,1], $\mu_f = 0$, $\sigma_f^2 = 1/3$ and thus

$$\rho_{\text{CEF, }exp}^{\star} = \frac{B^2 - K^2}{N/3 - K^2} \text{ (clipped)}.$$
(58)

If a deterministic worst-case guarantee is required (independent of the draw), note that $|\hat{u}_i| \leq 1$ almost surely, so the same bound as SA0/SA1 applies:

$$\rho_{\text{CEF, worst}}^{\star} = \frac{B^2 - K^2}{N} \text{ (clipped to } [0, 1]).$$
 (59)

E.5 **SUMMARY**

Here, we explain how to calculate the capacity of the proposed mechanism. Choose a gradient floor γ (e.g., f% of the arctangent peak so $\gamma = f \cdot \alpha_s/\pi$), compute $\delta(\gamma)$ and B via equation 51, measure $K = ||w||_2$, and then plug into the formula for the fault model of interest. If $\rho \leq \rho^*$, our mechanism keeps $|u_t| \leq \delta(\gamma)$ for all steps (deterministic case) or in expectation (stochastic case), thereby ensuring $\phi'(u_t) \geq \gamma$.

CONVERGENCE ANALYSIS WITH THE PROPOSED MECHANISM

We present that Stochastic Gradient Descent (SGD) and Gradient Descent (GD) optimizers derive gradients of fault-injected SNN models and induce the models to update weights when we adopt our mechanism to the SNN models.

F.1 SETUP AND NOTATION

We denote a data sample by (x,y), and the SNN by $f_W(\cdot)$ with parameters W. We include the (possibly stochastic) fragment + RMS transform T and analyze the expected objective $\tilde{L}(W) = \mathbb{E}_{(x,y),\,T}\big[\ell(f_W(T(x)),y)\big]$. At step t, with mini-batch estimator g_t and step size η_t , the update is $W_{t+1} = W_t - \eta_t g_t$. For spiking neurons, we write the pre-activation as $u_t := z_t - \vartheta$ and use a surrogate derivative $\phi'(u)$. Throughout, we assume the *gradient-corridor* condition $\phi'(u_t) \geq \gamma$ holds along the iterates, which is enforced by the fragment RMS bound and the per-layer effective weight norms $\|\hat{w}^{(\ell)}\|_2 \leq B^{(\ell)}$ referring to Appendix E. Symbols L, σ^2, μ are in Appendix C, D, and E, referenced only when required by a lemma or theorem.

F.2 BAND CONDITION ENFORCED BY OUR MECHANISM

With dynamic fragmentation and per-fragment RMS normalization $\|\tilde{x}_t\|_2 = \alpha_n$, each step satisfies

$$|u_t| = |\hat{w}^\top \tilde{x}_t + b - \vartheta| \le ||\hat{w}||_2 \alpha_n + m, \qquad m := |b - \vartheta|.$$
 (60)

Defining

$$B := \frac{\delta(\gamma) - m}{\alpha_n} \quad \text{(requires } \delta(\gamma) > m), \tag{61}$$

One obtains the following corridor-invariance lemma.

Lemma 4 (Corridor invariance). If $\|\hat{w}\|_2 \leq B$, then $|u_t| \leq \delta(\gamma)$ for all fragments t, hence $\phi'(u_t) \geq \gamma$. Proof. Combine equation 60 with equation 61 and the definition of $\delta(\gamma)$ in equation 17.

F.3 OPTIMIZATION OBJECTIVE AND ASSUMPTIONS

Let $\tilde{\mathcal{L}}(W) := \mathbb{E}_{(x,y),\mathcal{T}}[\ell(f_W(\mathcal{T}(x)),y)]$, where \mathcal{T} denotes the (possibly stochastic, data/model-aware) transformation induced by our mechanism (e.g., masks and RMS scaling). Assume:

Assumption 2 (*L*-smoothness). $\nabla \mathcal{L}$ is *L*-Lipschitz.

Assumption 3 (Unbiased mini-batch gradients, bounded variance). $\mathbb{E}[g_t | W_t] = \nabla \tilde{\mathcal{L}}(W_t)$ and $\mathbb{E}[\|g_t - \nabla \tilde{\mathcal{L}}(W_t)\|^2 | W_t] \leq \sigma^2$.

Assumption 4 (Corridor stability). For each layer, capacity constraints on the fault ratio ensure $\|\hat{w}^{(\ell)}\|_2 \leq B^{(\ell)}$, so Lemma holds layer-wise and $\phi'(u_t^{(\ell)}) \geq \gamma$ during training.

F.4 DESCENT LEMMA AND MASTER INEQUALITY

By L-smoothness and $W_{t+1} = W_t - \eta_t g_t$,

$$\tilde{\mathcal{L}}(W_{t+1}) \leq \tilde{\mathcal{L}}(W_t) - \eta_t \left\langle \nabla \tilde{\mathcal{L}}(W_t), g_t \right\rangle + \frac{L\eta_t^2}{2} \|g_t\|^2. \tag{62}$$

Taking expectation and using Assumption 2 with $\mathbb{E}\|g_t\|^2 = \|\nabla \tilde{\mathcal{L}}(W_t)\|^2 + \mathbb{E}\|g_t - \nabla \tilde{\mathcal{L}}(W_t)\|^2 \le \|\nabla \tilde{\mathcal{L}}(W_t)\|^2 + \sigma^2$ gives(Nesterov, 2014)

$$\mathbb{E}\big[\tilde{\mathcal{L}}(W_{t+1})\big] \leq \mathbb{E}\big[\tilde{\mathcal{L}}(W_t)\big] - \eta_t \left(1 - \frac{L\eta_t}{2}\right) \mathbb{E}\big[\|\nabla \tilde{\mathcal{L}}(W_t)\|^2\big] + \frac{L\eta_t^2}{2} \sigma^2.$$
 (63)

Theorem 2 (SGD convergence to stationarity). If $\eta_t \equiv \eta \in (0, 1/L]$, summing equation 63 over $t = 0, \dots, T-1$ yields

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\left[\|\nabla \tilde{\mathcal{L}}(W_t)\|^2 \right] \leq \frac{2\left(\tilde{\mathcal{L}}(W_0) - \tilde{\mathcal{L}}^{\star}\right)}{\eta T} + L \eta \sigma^2.$$
 (64)

With a Robbins-Monro schedule $(\sum_t \eta_t = \infty, \sum_t \eta_t^2 < \infty)$ we obtain $\lim_{T\to\infty} \min_{t< T} \mathbb{E} \|\nabla \tilde{\mathcal{L}}(W_t)\|^2 = 0$. Role of Assumption 3: by preventing artificial saturation $(\phi'(u) \approx 0)$, the corridor ensures that gradient signals remain informative until genuine stationarity (Bottou et al., 2018).

Theorem 3 (Monotone decrease for full-batch GD). In the deterministic case ($\sigma = 0$) with $\eta \in (0, 1/L]$,

$$\tilde{\mathcal{L}}(W_{t+1}) \leq \tilde{\mathcal{L}}(W_t) - \frac{\eta}{2} \|\nabla \tilde{\mathcal{L}}(W_t)\|^2, \tag{65}$$

so $\sum_t \|\nabla \tilde{\mathcal{L}}(W_t)\|^2 < \infty$ and every limit point of $\{W_t\}$ is stationary.

Corollary 1 (Linear rate under PL). If $\tilde{\mathcal{L}}$ satisfies the Polyak-Łojasiewicz (PL) inequality $\frac{1}{2}\|\nabla \tilde{\mathcal{L}}(W)\|^2 \geq \mu(\tilde{\mathcal{L}}(W) - \tilde{\mathcal{L}}^\star)$ for some $\mu > 0$ on the corridor-stable region, then for GD with $\eta \in (0, 1/L]$ (Karimi et al., 2020),

$$\tilde{\mathcal{L}}(W_t) - \tilde{\mathcal{L}}^* \le (1 - \eta \mu)^t \left(\tilde{\mathcal{L}}(W_0) - \tilde{\mathcal{L}}^* \right). \tag{66}$$

F.5 SUMMARY

Our mechanism enforces equation 17–equation 61 so that surrogate gradients do not vanish spuriously; Under standard smoothness/stochasticity assumptions, SGD converges to stationarity and GD decreases monotonically, with linear rates under PL. Capacity bounds on the fault ratio provide concrete regimes where the corridor assumption holds layer-wise (Neftci et al., 2019).