

A²SG: ADAPTIVE AND ASYMMETRIC SURROGATE GRADIENTS FOR TRAINING DEEP SPIKING NEURAL NETWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

Training deep spiking neural networks (SNNs) remains challenging due to sharp loss landscapes and temporal inconsistency caused by surrogate gradients. To address these challenges, we propose a unified framework: adaptive and asymmetric surrogate gradients (*A²SG*). The adaptive gradients adjust an effective window for spatio-temporal adaptation, reducing spatial gradient variation and maintaining directional consistency of gradients over time. The asymmetric gradients reflect neuronal dynamics by assigning larger gradients to neurons with higher membrane potentials, and we prove that they yield lower variation than symmetric surrogates. Our analysis further establishes a direct connection between local gradient variation and the curvature of the loss landscape, providing a principled explanation for how *A²SG* promotes convergence to flatter minima and improves generalization. We conduct extensive experiments on diverse models, including CNN-based and Transformer-based SNNs, across various tasks such as image classification using both static and neuromorphic datasets, as well as segmentation. The results demonstrate that *A²SG* consistently improves accuracy and energy efficiency, establishing it as a general and reliable solution for training deep SNNs.

1 INTRODUCTION

Spiking neural networks (SNNs) have emerged as energy-efficient next-generation neural networks that operate based on spikes (Maass, 1997). In particular, by leveraging the low-power characteristics of SNNs and the superior learning capabilities of deep neural networks (DNNs), deep SNNs have shown the potential for energy-efficient artificial intelligence in various fields (Park et al., 2020; Kim et al., 2020b; 2022b; Yao et al., 2025). Recently, deep SNNs have been successfully applied to various applications, including image segmentation (Kim et al., 2022b; Lei et al., 2025), object detection (Kim et al., 2020b; Su et al., 2023), and language modeling (Bal & Sengupta, 2024; Xing et al., 2024), as well as to diverse model architectures, such as Transformers (Zhou et al., 2023; Yao et al., 2025). These rapid advancements have been largely driven by the adoption of gradient-based training with surrogate gradients (Wu et al., 2018; Neftci et al., 2019).

Despite their essential role in training deep SNNs, research on effective surrogate gradient functions remains limited, contributing to the performance gap between DNNs and deep SNNs. Several studies have attempted to mitigate the mismatch between surrogate and true gradients by adaptively adjusting surrogate functions. However, existing approaches have predominantly focused on gradient sparsity for adaptation (Lian et al., 2023; Lin et al., 2023) or impose substantial computational overhead (Li et al., 2021), restricting the training performance or hindering practical deployment. Furthermore, few studies have designed surrogate functions that take into account the impact of surrogate gradients on generalization performance.

In this work, we introduce adaptive and asymmetric surrogate gradients (*A²SG*) to enhance the training of deep SNNs. The adaptive component leverages spatio-temporal adaptation, dynamically adjusting the surrogate gradient window to suppress spatial fluctuations of gradients and align their directions across timesteps. The asymmetric component allocates larger gradients to neurons with greater membrane potential, effectively prioritizing those closer to firing and promoting convergence to flatter minima. These designs are motivated by our theoretical analysis, which shows that a

054 larger variation in local gradient leads to sharper loss landscapes. Moreover, we demonstrate that the
 055 proposed asymmetric surrogate gradient exhibits lower gradient variation compared to its symmetric
 056 counterparts. In addition, we highlight temporal model collapse, which is caused by misaligned
 057 gradients across timesteps, as one of the obstacles for stable learning, motivating the need for spatio-
 058 temporal adaptation. By combining these, we establish a unified strategy that stabilizes optimization,
 059 promotes convergence to flatter minima, and improves generalization. We validate our proposed
 060 approaches through extensive experiments on both static and neuromorphic datasets, spanning convolutional
 061 neural networks (CNNs), Transformer-based models, and segmentation tasks. Across
 062 all benchmarks, A^2SG achieves consistent gains in accuracy and energy efficiency, highlighting its
 063 effectiveness as a general solution for reliable deep SNN training.

064 2 RELATED WORK

065 2.1 TRAINING DEEP SPIKING NEURAL NETWORKS

066 Recent studies have introduced deep SNN achieving both high performance and energy efficiency (Tavanaei et al., 2019). In these architectures, leaky integrate-and-fire (LIF) neurons are widely used for their computational simplicity and biological plausibility (Eqs. A1- A3). Deep SNNs have been applied to various tasks such as image classification (Hu et al., 2021; Fang et al., 2021a), object detection (Kim et al., 2020a;b), semantic segmentation (Kim et al., 2022b; Lei et al., 2025), and Transformer-based models (Zhou et al., 2023; Yao et al., 2025). Recent studies have adopted direct training based on spatio-temporal backpropagation (STBP) with surrogate gradients (Wu et al., 2018), as given in Eqs. A4 and A5. This method enables efficient training with fewer time steps, yet a performance gap remains compared to conventional DNNs. To mitigate this gap, several studies have focused on addressing the gradient mismatch problem arising from the adoption of surrogate gradients (Li et al., 2021; Lian et al., 2023). However, studies on gradient consistency during training have remained relatively limited. Especially, in STBP, parameter updates are obtained by aggregating gradient contributions from all timesteps, and inconsistent temporal gradients can generate conflicting signals, a phenomenon we term *temporal gradient confusion*. This inconsistency hinders stable optimization and degrades learning performance, highlighting the need for strategies that explicitly mitigate it.

084 2.2 SURROGATE GRADIENTS IN SPATIO-TEMPORAL BACKPROPAGATION (STBP)

085 Surrogate gradients have been employed to address the non-differentiability of spiking function (Eq. A2) during error backpropagation. Although the adoption of surrogate gradients has dramatically improved the performance of deep SNNs, their use remains limited by inconsistencies with the true gradients. To improve the learning performance, several studies have focused on adjusting the distribution of the membrane potential during training (Guo et al., 2022; 2023a;b; Zhao et al., 2025). Various regularization strategies have been employed, such as maximizing the information within the membrane potential (Guo et al., 2022) or minimizing quantization errors induced by the spike function (Guo et al., 2023a). In addition, batch normalization (Guo et al., 2023b) and KL loss (Zhao et al., 2025) were applied to mitigate the inter-batch and temporal discrepancies in the membrane potential distribution, respectively. While these studies improved performance by adjusting the membrane potential distribution, they did not fundamentally address the performance degradation inherent to surrogate gradients, highlighting the essential need for advancements in surrogate gradient design.

099 2.3 IMPROVING SURROGATE GRADIENTS DESIGN FOR DEEP SNNs

100 Most surrogate gradients adopt static and symmetric function shapes to approximate the Dirac delta, which represents the derivative of the spiking function ($\frac{\partial s[t]}{\partial u[t]}$). These functions preserve a constant area within the effective window $[V_{th} - \beta, V_{th} + \beta]$, with representative examples being the box-car (*BOX*) and triangle (*TRI*) functions (Eqs. A6-A7.) To mitigate vanishing gradients, Guo et al. (2024) proposed directly delivering gradients to shallow layers, though gradient mismatch remains unresolved. Beyond static functions, adaptive strategies have been explored to improve training further (Li et al., 2021; Lian et al., 2023; Lin et al., 2023; jia). Dspike (Li et al., 2021) employs finite-difference gradients to align surrogates with true gradients via cosine similarity, though its

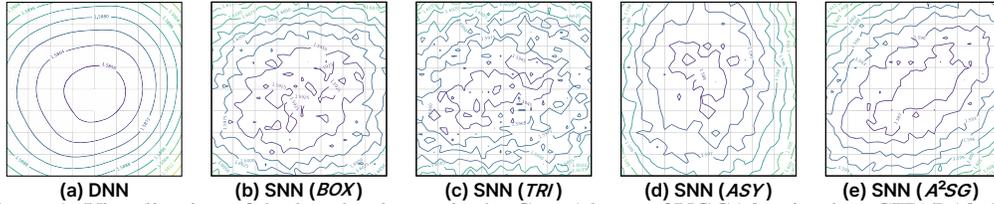


Figure 1: Visualization of the loss landscape in the Conv1 layer of VGG16 trained on CIFAR10. The results demonstrate that our adaptive surrogate gradient leads to a flatter loss landscape compared to conventional surrogate gradients.

high computational cost limits scalability. To reduce this overhead, other studies propose adjusting the effective window width to regulate gradient sparsity during training (Lian et al., 2023; Lin et al., 2023; jia). However, most of these approaches have focused on sparsity control to avoid gradient vanishing or explosion, relying on sparsity as an indirect indicator of learning quality. Moreover, existing efforts have primarily concentrated on symmetric functions that mimic differential operators, while the impact of function shape itself remains underexplored.

2.4 FLAT MINIMA AND GENERALIZATION

Although gradient-based direct learning has significantly improved the performance, generalization remains a critical challenge for deep SNNs. Numerous studies have reported that models converging to flat minima tend to achieve better generalization (Hochreiter & Schmidhuber, 1997; Keskar et al., 2017; Chaudhari et al., 2019). Flatness is commonly assessed via the Hessian spectrum (Ghorbani et al., 2019), but direct Hessian computations are intractable for large models. As a practical alternative, the Fisher information matrix (FIM) (Eq. A8) is frequently used, since its eigenvalues are known to capture the curvature of the loss landscape, with smaller values indicating flatter minima (Liao et al., 2018; Karakida et al., 2019; Martens, 2020; Kim et al., 2022a). In addition to measurement, several training procedures aim to encourage convergence to flat minima. These include entropy-based biasing toward wide valleys (Chaudhari et al., 2019) and curvature-aware updates using FIM-based criteria (Kim et al., 2022a). Related studies have also linked sharp minima to large-batch training and poor generalization (Keskar et al., 2017), further motivating flatness-oriented training strategies. Overall, prior works suggest that guiding optimization toward flat minima is an effective approach for improving generalization.

3 SHARP LOSS LANDSCAPE FROM SURROGATE GRADIENT LEARNING

Deep SNNs trained with surrogate gradients tend to converge to sharper loss landscapes than DNNs. The flatness of the loss landscape is characterized by its curvature, quantified by the Hessian $\mathbf{H} = \nabla_{\mathbf{w}}^2 L(\mathbf{w})$, with respect to the parameter vector \mathbf{w} . For clarity, we analyze the second derivative of the loss with respect to a single weight, which corresponds to a diagonal entry of the Hessian. By the chain rule, the second derivative of the loss with respect to a weight w can be described as

$$\frac{\partial^2 L}{\partial w^2} = \frac{\partial^2 L}{\partial \phi^2} (\phi'(u)x)^2 + \frac{\partial L}{\partial \phi} \phi''(u)x^2, \quad (1)$$

where x denotes the pre-synaptic activation and $u = wx$. For DNNs with an activation function $\phi(u)$, assume the first and second derivatives are bounded, i.e., $|\phi'(u)| \leq c_1$ and $|\phi''(u)| \leq c_2$ for finite constants c_1 and c_2 , which yields

$$\left| \frac{\partial^2 L}{\partial w^2} \right|_{DNN} = \mathcal{O}(x^2). \quad (2)$$

For deep SNNs, a surrogate gradient function $f(u)$ introduces an effective window of width β to approximate the non-differentiable spike function. As shown in Sec. A.5, any symmetric surrogate with fixed area satisfies

$$\|H'\|_{\infty} = \|f\|_{\infty} = \Omega(\beta^{-1}) \quad \text{and} \quad \|H''\|_{\infty} = \|f'\|_{\infty} = \Omega(\beta^{-2}). \quad (3)$$

Substituting them into the chain-rule expansion (Eq. 1) yields

$$\left| \frac{\partial^2 L}{\partial w^2} \right|_{SNN} = \Omega\left(\frac{x^2}{\beta^2}\right). \quad (4)$$

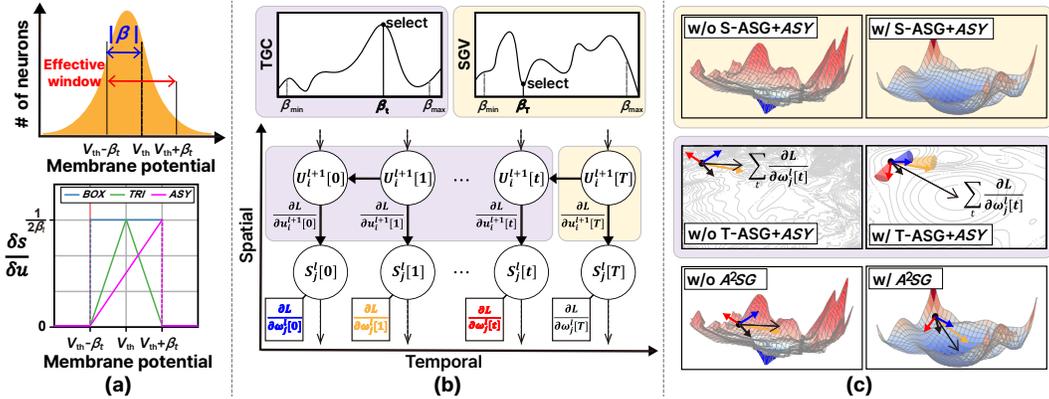


Figure 2: Overview of the proposed A^2SG framework. (a) Effective window (red) is modulated by the parameter β (blue); three shapes (*BOX*, *TRI*, and *ASY*) are shown. (b) Temporal adaptive surrogate gradient (T-ASG) selects β_t that maximizes temporal gradient consistency (TGC). While Spatial adaptive surrogate gradient (S-ASG) selects β_t that minimizes spatial gradient variation (SGV). (c) S-ASG+ASY promotes flat minima, T-ASG+ASY stabilizes gradient directions, and A^2SG achieves robust convergence by combining both.

Most prior works adopt a *narrow effective window* ($\beta < 1$) to better approximate the Dirac delta around the threshold, which empirically improves training stability and convergence (Wu et al., 2018; Neftci et al., 2019). Under this condition, Eq. 4 shows that the Hessian magnitude is amplified by a factor of $1/\beta^2$ in surrogate-trained SNNs, indicating a sharper loss landscape than DNNs. In addition, the binary and temporally sparse nature of spikes concentrates gradients and increases their variation, further sharpening the landscape, as shown in Fig. 1. Consistent with this analysis, *TRI* yields a sharper loss landscape than *BOX* because, under area normalization, its steeper slopes imply larger curvature. This is derived in Sec. A.6 and confirmed empirically through Fig. 1-(b) and (c).

4 A^2SG : ADAPTIVE AND ASYMMETRIC SURROGATE GRADIENTS

We introduce A^2SG to address two limitations in surrogate-based SNN training: sharp loss landscapes and temporal gradient confusion. The adaptive component consists of two policies: spatial and temporal adaptations to address sharp loss landscapes and temporal gradient confusion, respectively. The spatial adaptation adjusts the surrogate effective window to reduce the variance of the local gradient ($\frac{\partial L}{\partial u}$), encouraging convergence to flatter minima. The temporal adaptation aligns per-timestep local gradients to mitigate temporal inconsistency. We demonstrate that the dispersion of local gradients affects the curvature of the loss landscape and that temporal gradients are formed by aggregating these local gradients. Thus, adjusting β during training effectively reduces spatial variability and enhances temporal alignment, improving generalization under non-stationary dynamics.

In addition, we analyze how the functional shape of the surrogate affects local gradient variation and introduce an asymmetric surrogate that considers neuronal dynamics. By allocating larger gradients to neurons with greater accumulated membrane potential, the asymmetric form further suppresses the gradient variance and alleviates sharpness. When utilized together, the spatio-temporal adaptation and the asymmetric components complement each other, effectively addressing the significant limitations of surrogate-gradient learning. They stabilize the optimization process and yield flatter solutions with improved generalization. The overall framework of A^2SG is illustrated in Fig. 2.

4.1 RELATION BETWEEN VARIATION OF LOCAL GRADIENT AND FLATNESS OF LOSS LANDSCAPE

To analyze how the variability of local gradients affects the flatness of the loss landscape, we focus on the coefficient of variation (CV) of the local gradients. For notational simplicity, we consider a fully connected (FC) layer, but the principle naturally extends to other neural network layers and architectures. Let $W \in \mathbb{R}^{m \times n}$ denote the weight matrix of the FC layer. The gradient with respect to W , vectorized, is given by: $\mathbf{g} = \text{vec}\left(\frac{\partial L}{\partial W}\right) = \mathbf{a}_{\text{in}} \otimes \boldsymbol{\delta}$, where $\mathbf{a}_{\text{in}} \in \mathbb{R}^n$ is the input vector and

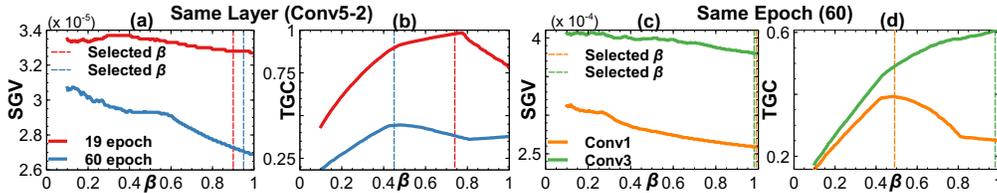


Figure 3: Graphs of SGV and TGC as a function of β . (a) and (b) represent SGV and TGC, at different epochs for the Conv5-2 layer, while (c) and (d) compare Conv1 and Conv3 at epoch 60. Dashed lines indicate the selected β through the proposed adaptive method in each training case.

$\delta \in \mathbb{R}^m$ is the backpropagated error vector. The FIM is then defined as:

$$\mathbf{F} = \mathbb{E}[\mathbf{g}\mathbf{g}^\top] = \mathbb{E}[(\mathbf{a}_{\text{in}} \otimes \delta)(\mathbf{a}_{\text{in}} \otimes \delta)^\top]. \quad (5)$$

For analytical clarity, the δ can be decomposed into its mean and a zero-mean fluctuation term: $\delta = \mu\mathbf{1} + \epsilon$, where the μ is the mean of δ , $\mathbf{1}$ is the all-ones vector and ϵ denotes a small perturbation. Substituting this into the definition of \mathbf{F} , we obtain:

$$\mathbf{F} = \mu^2 \mathbb{E}[(\mathbf{a}_{\text{in}} \otimes \mathbf{1})(\mathbf{a}_{\text{in}} \otimes \mathbf{1})^\top] + \mathbf{R} = \mathbf{F}_0 + \mathbf{R}, \quad (6)$$

where \mathbf{F}_0 is a rank-1 matrix and the perturbation \mathbf{R} is bounded as $\|\mathbf{R}\|_2 \leq c\mu\text{CV}(\delta)$ for some constant c . This shows that as $\text{CV}(\delta)$ becomes smaller, the FIM approaches the rank-1 matrix \mathbf{F}_0 , indicating that the loss landscape has a dominant curvature direction. By matrix perturbation theory (Greenbaum et al., 2020), the largest eigenvalue of the FIM is bounded as:

$$\lambda_{\max}(F) \leq \mu^2 \lambda_{\max}(\mathbb{E}[(\mathbf{a}_{\text{in}} \otimes \mathbf{1})(\mathbf{a}_{\text{in}} \otimes \mathbf{1})^\top]) + c\mu^2\text{CV}(\delta), \quad (7)$$

where $\lambda_{\max}(\cdot)$ denotes the largest eigenvalue operator. Therefore, the largest eigenvalue grows linearly with $\text{CV}(\delta)$, and reducing the CV of local gradients directly leads to a flatter loss landscape.

4.2 SPATIO-TEMPORAL ADAPTIVE SURROGATE GRADIENTS (ST-ASG)

As discussed in the previous section, reducing local gradient variability alleviates the sharpness of the loss landscape. In addition, temporal gradient confusion can be mitigated by promoting alignment of local gradients across timesteps. To achieve this, we introduce two metrics: spatial gradient variation (SGV) and temporal gradient consistency (TGC), which guide spatial and temporal adaptation, respectively. SGV is defined as follows:

$$\text{SGV}^{(l)}[T] := \frac{\text{Var}(\delta^l[T])}{\text{Mean}(|\delta^l[T]|)}, \quad (8)$$

where $\delta^{(l)}[T]$ denotes the backpropagated error in layer l at the last timestep T . To improve computational efficiency in practice, SGV employs the variance rather than the standard deviation in the denominator, differing from the conventional CV. On the other hand, TGC at timestep t is defined as the cosine similarity between the local gradients at adjacent timesteps:

$$\text{TGC}^{(l)}[t] := \cos(\delta^{(l)}[t], \delta^{(l)}[t+1]), t \in [1, T-1]. \quad (9)$$

With these definitions, spatial adaptation is applied using SGV to suppress local gradient variation at the last timestep, where activations and gradients are relatively stable. In this case, the adaptation objective is to minimize SGV. Conversely, temporal adaptation is guided by TGC, which promotes alignment of local gradients between adjacent timesteps. In this context, the goal of adaptation is to maximize TGC. Since spatial adaptation is applied at the last timestep, it provides a stable reference direction for the temporal gradients. Temporal adaptation, applied to preceding timesteps ($t < T$), then aligns the local gradients with this reference. In this way, spatio-temporal adaptation is achieved by anchoring the global gradient trajectory to the stable direction obtained at the last timestep, while simultaneously enforcing temporal consistency across earlier timesteps.

To realize this spatio-temporal adaptation in practice, we propose to adjust the width (β) of the effective window. Our method is motivated by the observation that both SGV and TGC can be expressed as functions of β . As illustrated in Fig. 3, these functions vary unpredictably with the training dynamics: the same layer exhibits different functional shapes across epochs, and even within a single epoch, distinct patterns may emerge across layers. To robustly identify suitable values of β under such variability, we employ a Bayesian search strategy. Implementation details of the adaptive method and the Bayesian optimization procedure are provided in Secs. A.7 and A.8, respectively.

4.3 ASYMMETRIC SURROGATE GRADIENTS

Symmetric functions, such as *TRI* and *BOX*, have been widely used as surrogate gradients in STBP, providing gradients based solely on the distance between a neuron’s membrane potential and the threshold. However, they fail to account for neuronal dynamics such as integration and firing. Thus, the relative magnitude of the accumulated membrane potential is not effectively reflected in the training process. To solve this problem, we propose an asymmetric (*ASY*) surrogate, defined as

$$\frac{\partial s}{\partial u} = f(u, \beta) = \frac{1}{2\beta} \cdot (u - V_{\text{th}}) + h, \quad u \in [V_{\text{th}} - \beta, V_{\text{th}} + \beta]. \quad (10)$$

h is empirically determined considering the gradient sparsity of each model. The proposed *ASY* function produces a larger gradient when the membrane potential is more highly accumulated, allowing the learning algorithm to capture each neuron’s contribution based on its dynamic behavior.

This neuron’s behavior-aware design also leads to a reduction in gradient variability. By concentrating gradient values in regions where membrane potential is high (and spike likelihood is greater), the *ASY* function avoids spreading gradients across irrelevant low-activity regions. This focused gradient allocation reduces unnecessary variability, resulting in more stable training.

This intuition is formalized in the following theoretical results.

Theorem 1 (CV-Minimizing Symmetric Function under Area and Boundary Constraints). *Let $f : [a, b] \rightarrow \mathbb{R}_{\geq 0}$ be a function. Let $f_{\text{asy}}(u)$ and $f_{\text{sym}}(u)$ be asymmetric and symmetric surrogate gradient functions defined over $[a, b]$, satisfying the boundary condition $f(a) = f(b) = 0$, non-negativity $f(u) \geq 0$, and area constraint $\int_a^b f(u) du = c$, where $a = \theta - \beta$ and $b = \theta + \beta$. Suppose the membrane potential $u \sim \mathcal{N}(\mu, \sigma^2)$ with $\mu < a$, so that $p(u)$ is decreasing on $[a, b]$. Then the unique function f^* that minimizes the CV over all such admissible functions is the symmetric triangular function.*

Proof. Please refer to Thm. A1. □

Thm. 1 shows that, under area and boundary constraints, the symmetric function that minimizes gradient CV is a triangular function. This sets a lower bound of CV for symmetric functions under the given constraints. Based on this fact, we verify that asymmetric has a lower CV than symmetric.

Theorem 2 (CV Comparison of Asymmetric and Symmetric Surrogates). *Let $f_{\text{asy}}(u)$ and $f_{\text{sym}}(u)$ be asymmetric and symmetric surrogate gradient functions defined over $[a, b]$ under the same constraints as in Thm. 1. Then, under a linear approximation of the Gaussian, where $L = b - a$ and $\kappa = a - \mu$, we have:*

$$CV_{\text{asy}} < CV_{\text{sym}} \quad \text{if } L\kappa > \sigma^2.$$

Proof. Please refer to Thm. A2. □

Thms. 1 and 2 demonstrate that the proposed asymmetric surrogate gradient, by incorporating membrane potential accumulation, achieves a lower CV than its symmetric counterparts. These theoretical findings indicate that the asymmetric surrogate gradient not only better reflects neuronal dynamics but also promotes more stable and efficient learning, thereby facilitating convergence to flatter minima. Experimental validations of Thms. 1 and 2 are provided in Fig. 5-(a), where *ASY* function consistently exhibits lower gradient variance than *TRI* function. Fig. 4 presents $L\kappa$ and σ^2 of Conv1 and Conv5-2 during training on VGG16 with CIFAR10. The graphs show that the condition $L\kappa > \sigma^2$ in Thm. 2 becomes satisfied across layers as training progresses. This confirms that our assumption is supported by the experimental results.

5 EXPERIMENTS

We evaluated the effectiveness of the proposed method on various datasets, including static image datasets such as CIFAR10, CIFAR100 (Krizhevsky et al., 2009), and ImageNet (Deng et al., 2009) as

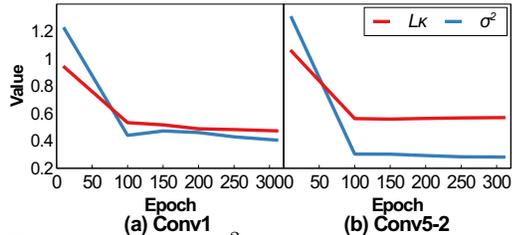


Figure 4: $L\kappa$ and σ^2 values across epochs. (a) and (b) correspond to Conv1 and Conv5-2, respectively.

Figure 4 presents $L\kappa$ and σ^2 of Conv1 and Conv5-2 during training on VGG16 with CIFAR10. The graphs show that the condition $L\kappa > \sigma^2$ in Thm. 2 becomes satisfied across layers as training progresses. This confirms that our assumption is supported by the experimental results.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

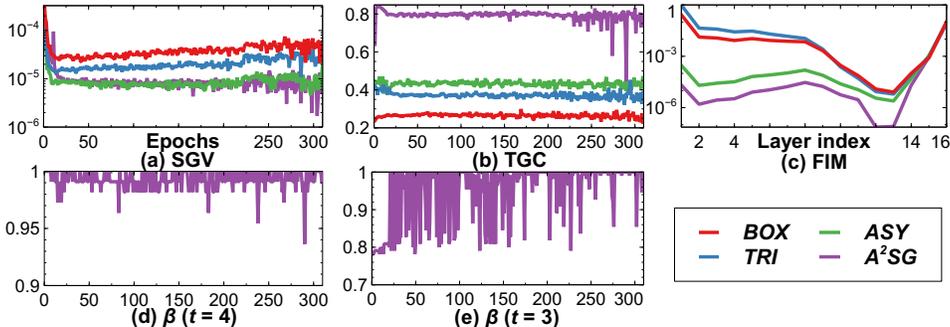


Figure 5: Comparison of SGV, TGC, FIM, and β dynamics across different surrogate gradient functions at Conv5-2 layer. (a) SGV over epochs, (b) TGC over epochs, (c) FIM across layer indices, (d) β dynamics at $t = 4$, and (e) β dynamics at $t = 3$.

well as a neuromorphic dataset such as CIFAR10-DVS (Li et al., 2017). We conducted experiments on both CNN and Transformer models. To further show the versatility of our method, we also evaluated it on the ADE20K (Zhou et al., 2017) dataset for semantic segmentation. For more details about the experimental setup, please refer to Secs. A.9 and A.10.

5.1 EFFECT OF A^2SG ON GRADIENT DYNAMICS AND FEATURE LEARNING

We analyze the effect of A^2SG on gradient variation and temporal consistency. As shown in Fig. 5-(a) and (b), A^2SG maintains low SGV and high TGC throughout training. Furthermore, in Fig. 5-(c), measurement of the maximum eigenvalue of the FIM for each layer reveals that A^2SG consistently achieves the lowest eigenvalues throughout the network. These results corroborate our theoretical findings, demonstrating that lower CV directly leads to convergence toward flatter minima in the loss landscape. In addition, Fig. 5-(d) and (e) illustrate the adaptive selection of β across training epochs, where β is chosen to minimize SGV and maximize TGC, respectively. Additional analysis of the ASY function and the results for the Conv1 and Conv3 layers are provided in Secs. A.12 and A.13, respectively.

To further examine how improved gradient dynamics translate into feature representations, we perform a t-SNE visualization of the learned feature representations (Fig. 6). Models trained with A^2SG exhibit class-separable features at early timestep $t = 1$ and by $t = 4$ the features of each class consolidate into compact and well-separated clusters. These properties indicate that gradients are effectively propagated to deeper layers and maintain coherent directions across timesteps, facilitating convergence to flat minima and improving training ability. In contrast, BOX exhibits significant class overlap and less clear separation. This demonstrates that A^2SG enhances both training stability and generalization performance, as visualized by more robust and discriminative feature representations.

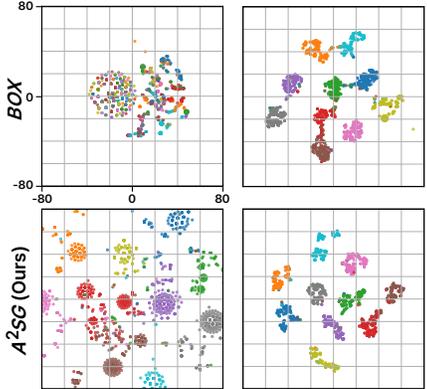


Figure 6: Comparison of t-SNE for VGG16 on CIFAR10 using BOX (top) and A^2SG (bottom) at $t = 1$ and $t = 4$.

5.2 COMPARISON WITH OTHER ADAPTIVE SURROGATE GRADIENTS

Tab. 1 compares the accuracy of various adaptive surrogate gradient methods on CIFAR10 and CIFAR100 with ResNet18 and ResNet19. When applied to BOX and TRI , our ST-ASG consistently outperforms previous adaptive surrogate gradient methods on both datasets. This advantage arises from its spatio-temporal adaptation. The spatial policy reduces variability in local gradients, encouraging convergence toward flatter minima. The temporal policy aligns gradient directions across timesteps, stabilizing the global update. Furthermore, A^2SG achieves the highest accuracy among all methods. On ResNet19, our approach outperforms LSG by 1.57% on CIFAR10 (96.74% vs. 95.17%) and by 4.20% on CIFAR100 (81.05% vs. 76.85%). These results demonstrate the strong generalization capability of our method over existing adaptive surrogate gradient approaches.

Table 1: Comparison with other adaptive surrogate gradient methods (ARC: Arctangent, ST-ASG: Our spatio-temporal adaptive surrogate gradient).

| Datasets | Architectures | Functions | Methods | Timesteps | Acc. (%) |
|----------|---------------|-----------|--------------------------|-----------|-------------------|
| CIFAR10 | ResNet18 | ARC | Dspike (Li et al., 2021) | 4 | 93.66±0.05 |
| | | | CPNG (Lin et al., 2023) | 6 | 94.10±0.05 |
| | ResNet19 | TRI | ST-ASG | 4 | 96.41±0.16 |
| | | | LSG (Lian et al., 2023) | 4 | 95.17±0.05 |
| | | | ST-ASG | 4 | 96.01±0.05 |
| | | | A²SG | 4 | 96.74±0.05 |
| CIFAR100 | ResNet18 | ARC | Dspike (Li et al., 2021) | 4 | 73.35±0.14 |
| | | | CPNG (Lin et al., 2023) | 6 | 75.37±0.05 |
| | ResNet19 | TRI | ST-ASG | 4 | 80.46±0.06 |
| | | | LSG (Lian et al., 2023) | 4 | 76.85±0.10 |
| | | | ST-ASG | 4 | 78.60±0.16 |
| | | | A²SG | 4 | 81.05±0.05 |

Table 2: Comparison with current state-of-the-art approaches on CIFAR10/100.

| Datasets | Architectures | Methods | Timesteps | Acc. (%) |
|------------------------|---------------|-------------------------------|-----------|-------------------|
| CIFAR10 | VGG16 | IM (Guo et al., 2022) | 5 | 93.85 |
| | | RMP (Guo et al., 2023a) | 4 | 93.33 |
| | | MPBN (Guo et al., 2023b) | 4 | 94.44 |
| | | A²SG | 4 | 95.29±0.05 |
| | | IM (Guo et al., 2022) | 4 | 95.40 |
| | ResNet19 | RMP (Guo et al., 2023a) | 4 | 95.51 |
| | | TET (Deng et al., 2022) | 4 | 94.44 |
| | | TAB (Jiang et al., 2024) | 4 | 94.76 |
| | | ShortcutBP (Guo et al., 2024) | 2 | 95.36 |
| | | MPD-AGL (Jiang et al., 2025) | 2 | 96.18 |
| | | | 4 | 96.35 |
| | | | 6 | 96.54 |
| | | A²SG | 2 | 96.34±0.02 |
| | | A²SG | 4 | 96.74±0.05 |
| CIFAR100 | VGG16 | IM (Guo et al., 2022) | 5 | 70.18 |
| | | RMP (Guo et al., 2023a) | 4 | 72.55 |
| | | MPBN (Guo et al., 2023b) | 4 | 74.74 |
| | | A²SG | 4 | 75.21±0.08 |
| | ResNet19 | RMP (Guo et al., 2023a) | 4 | 78.28 |
| | | TET (Deng et al., 2022) | 4 | 74.47 |
| | | TAB (Jiang et al., 2024) | 4 | 76.81 |
| | | ShortcutBP (Guo et al., 2024) | 2 | 77.79 |
| | | MPD-AGL (Jiang et al., 2025) | 2 | 78.84 |
| | | | 4 | 79.72 |
| | 6 | 80.49 | | |
| A²SG | 2 | 79.18±0.01 | | |
| A²SG | 4 | 81.05±0.05 | | |

5.3 COMPARISON WITH STATE-OF-THE-ART METHODS

As reported in Tabs. 2 and 3, our method outperforms prior approaches on CIFAR10 and CIFAR100 in accuracy and achieves both higher accuracy and lower power consumption on ImageNet. In the case of E-SpikeFormer (Tab. 3), it adopted integer LIF neurons with multiple thresholds, which is distinct from the conventional SNNs with LIF neurons. To efficiently train on large-scale datasets, it uses integer values as a substitute for the temporal spike trains. Thus, this mechanism implicitly incorporates temporal dynamics, even when $T=1$. The model can be regarded as operating with an implicit time step equal to the maximum integer spike count D , rather than a single time step. Based on this perspective, we aligned the time step of the conventional LIF neurons with the integer activation value of I-LIF. For example, when $T \times D$ is 1x4 for I-LIF, we applied S-ASG to neurons with an activation value of four, corresponding to the last time step. We then sequentially applied T-ASG to neurons with activation values of three, two, and one. For the effective window (β), we set it to be centered on each threshold, as in LIF with a single threshold. For example, if th_i is a threshold at integer i , the effective window of th_i is set to $[th_i - \beta_i, th_i + \beta_i]$. From this state, we changed β_i through our adaptive method.

Tab. 4 shows that our method maintains higher accuracy on CIFAR10-DVS with only four timesteps. In Tab. 5, our approach also attains higher mIoU and reduced power consumption on ADE20K,

Table 3: Comparison with Transformer-based SNNs on ImageNet. Following E-Spikeformer Yao et al. (2025), timesteps are denoted as $T \times D$, where T is the number of timesteps and D indicates the upper bound of integer activations.

| Architecture | Methods | Param (M) | Power (mJ) | Time Steps | Acc. (%) |
|--------------|---|--------------|--------------|--------------|-------------------|
| Transformer | SpikFormer (Zhou et al., 2023) | 66.3 | 21.5 | 4×1 | 74.8 |
| | Meta-SpikeFormer (Yao et al., 2024) | 31.3 | 32.8 | 4×1 | 77.2 |
| | E-SpikeFormer(Yao et al., 2025) | 10.0 | 3.0 | 1×4 | 78.5 |
| | E-SpikeFormer(Yao et al., 2025) | 173.0 | 35.6 | 1×4 | 84.7 |
| | E-SpikeFormer + A^2SG | 10.0 | 2.78 | 1×4 | 78.61±0.01 |
| | E-SpikeFormer + A^2SG | 173.0 | 35.64 | 1×4 | 85.43 |

Table 4: Comparisons with other works on CIFAR10-DVS (* denotes our implementation).

| Datasets | Architectures | Methods | Timesteps | Acc. (%) |
|-------------|---------------|---------------------------------|-----------|-------------------|
| CIFAR10-DVS | VGGsNN | STBP-tdBN (Zheng et al., 2021)* | 4 | 81.30±1.00 |
| | | HSD (Zhong et al., 2024) | 5 | 81.10 |
| | | TMC (Yan et al., 2025) | 4 | 81.76 |
| | | A^2SG | 4 | 82.36±0.01 |

Table 5: Performance of segmentation on ADE20K. These methods use the pre-trained models on ImageNet as the backbone, then add segmentation heads for fine-tuning.

| Dataset | Methods | Param (M) | Power (mJ) | Time Steps | MIoU. (%) |
|---------|---|-----------|-------------|--------------|--------------|
| ADE20K | Meta-SpikeFormer (Yao et al., 2024) | 16.5 | 88.1 | 4×1 | 33.6 |
| | E-SpikeFormer (Yao et al., 2025) | 11.0 | 27.2 | 1×4 | 40.1 |
| | E-SpikeFormer + A^2SG | 11.0 | 25.2 | 1×4 | 40.94 |

demonstrating its effectiveness for segmentation as well as its energy efficiency. Moreover, Fig. A2 illustrates qualitative improvements in segmentation when applied to E-SpikeFormer. Notably, A^2SG converges rapidly, achieving competitive accuracy with only two timesteps, comparable to other state-of-the-art methods that require four timesteps. Overall, these results demonstrate that by improving the design of surrogate gradients, our method provides a general and principled strategy for enhancing the training performance of deep SNNs across diverse architectures and tasks.

5.4 ABLATION STUDIES

Tab. 6 summarizes ablation results on CIFAR10 with VGG16, highlighting the contributions of the adaptive and asymmetric components. Incorporating the spatial adaptive surrogate gradient (S-ASG) improves accuracy and reduces spike count, while the temporal adaptive surrogate gradient (T-ASG) alone preserves accuracy with a slight increase in spikes. Their combination (ST-ASG) further stabilizes training, and adding the asymmetric surrogate (A^2SG) yields the highest accuracy (95.29%) with the lowest spike count, confirming the benefit of integrating all components. To estimate the computational overhead of the proposed method, we measured the wall-clock time of one training epoch for each ablation case. As shown in the table, the proposed search method incurs a computational overhead of up to approximately 15% compared to the baseline.

Table 6: Ablation study on CIFAR10/100 with VGG16, comparing spatial (S-ASG), temporal (T-ASG), spatio-temporal adaptation (ST-ASG), and ST-ASG with ASY (A^2SG).

| | Methods | Acc. (%) | # of Spikes ($\times 10^3$) | Latency (sec/epoch) |
|----------|----------------------------------|-------------------|-------------------------------|---------------------|
| CIFAR10 | BOX (Baseline) | 94.84±0.05 | 94.6±1.0 | 74 (+0%) |
| | w/ S-ASG | 94.97±0.04 | 80.0±0.8 | 82 (+11%) |
| | w/ T-ASG | 94.94±0.05 | 98.7±1.8 | 83 (+12%) |
| | w/ ST-ASG | 94.98±0.03 | 93.6±2.3 | 85 (+15%) |
| | A^2SG (Ours) | 95.29±0.04 | 84.9±1.8 | 85 (+15%) |
| CIFAR100 | BOX (Baseline) | 74.24±0.08 | 100.3±1.2 | 74 (+0%) |
| | w/ S-ASG | 74.57±0.08 | 93.2±0.7 | 82 (+11%) |
| | w/ T-ASG | 74.54±0.06 | 104.9±2.2 | 83 (+12%) |
| | w/ ST-ASG | 74.73±0.07 | 100.2±0.2 | 85 (+15%) |
| | A^2SG (Ours) | 75.21±0.08 | 100.2±0.4 | 85 (+15%) |

5.5 NOISE ROBUSTNESS

In this section, we analyze the noise robustness of A^2SG . We experimented with deletion noise, removing a proportion of spikes from each layer, and report the results in Tab. A3.

486 Compared to the *BOX*, A^2SG achieves higher accuracy
 487 under the deletion noise, while also exhibiting a smaller
 488 reduction in total spike counts. In this section, we ana-
 489 lyze the noise robustness of A^2SG . We experimented with
 490 deletion noise, removing a proportion of spikes from each
 491 layer, and report the results in Tab. A3.

492 Fig. 7 illustrates the weight distributions, where A^2SG
 493 shows lower variance than *BOX*. This distribution indi-
 494 cates reduced reliance on specific weights, thereby im-
 495 proving generalization and enhancing robustness to errors (Tsai et al., 2021).

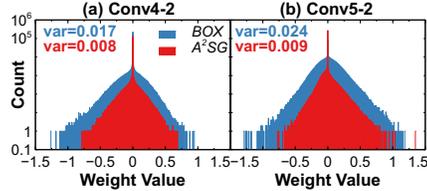


Figure 7: Weight distribution of *BOX* and A^2SG . (a) and (b) are Conv4-2 and Conv5-2 layers

5.6 COMPATIBILITY

498 To validate the compatibility of our approach, we additionally applied A^2SG to other neuron mod-
 499 els and learning methods. Specifically, we evaluated it with PLIF neurons (Fang et al., 2021b) or
 500 RMP-loss (Guo et al., 2023a), and the corresponding results are reported in Tab. 7 and Tab. 8, re-
 501 spectively. Experimental results of PLIF or RMP-loss show improvements in both accuracy and
 502 efficiency, which is consistent with other experimental results. These experiments demonstrate that
 503 our method is compatible with various neuron models and training methods to improve both accu-
 504 racy and efficiency.

Table 7: Comparison between LIF and PLIF on CIFAR10 with VGG16.

| Neurons | Methods | Acc. (%) | # of Spikes ($\times 10^3$) |
|--------------------------|----------------|----------------------------------|--------------------------------|
| LIF | Baseline | 94.84 \pm 0.05 | 94.6 \pm 1.0 |
| | A^2SG (Ours) | 95.29\pm0.04 | 84.9\pm1.8 |
| PLIF Fang et al. (2021b) | Baseline | 94.99 \pm 0.03 | 91.6 \pm 7.4 |
| | A^2SG (Ours) | 95.33\pm0.01 | 82.2\pm0.5 |

Table 8: Comparison between tdBN and RMP-loss on CIFAR100 with VGG16.

| Methods | Acc. (%) | # of Spikes ($\times 10^3$) |
|----------------------------------|----------------------------------|---------------------------------|
| tdBN | 74.24 \pm 0.08 | 100.3 \pm 1.2 |
| RMP Guo et al. (2023a) | 74.39 \pm 0.07 | 102.8 \pm 2.9 |
| tdBN + A^2SG | 75.21\pm0.08 | 100.2\pm0.4 |
| RMP + A^2SG | 75.25\pm0.03 | 102.3\pm1.1 |

5.7 SENSITIVITY ANALYSIS

514 In this section, we conduct a sensitivity anal-
 515 ysis of the hyperparameters of A^2SG . All re-
 516 sults reported in the experiments use the fol-
 517 lowing default configuration: β update fre-
 518 quency of 1 epoch, Bayesian optimization pa-
 519 rameters (n_{obs}, n_{eval}) set to (100, 150), and a
 520 search radius δ of 0.05. As shown in Tab. 9,
 521 increasing the update frequency to 50 or 100
 522 epochs widens the interval between β updates,
 523 which leads to accuracy drops within 0.2%.
 524 Similarly, when varying (n_{obs}, n_{eval}) to (10,
 525 15) and (300, 450), the accuracy difference re-
 526 mained within 0.2%. Finally, we examined the
 527 sensitivity of δ , which indicates the search width around β . The results obtained by varying it to
 528 0.01 and 0.1 are reported in Tab. 9. Overall, the experimental observations indicate that A^2SG
 529 demonstrates low sensitivity to hyperparameter variations.

Table 9: Comparison with different β update frequencies (epoch) and Bayesian optimization hyperparameter ($n_{obs}, n_{eval}, \delta$) settings on CIFAR10 with VGG16.

| Epoch | n_{obs} | n_{eval} | δ | Acc. (%) |
|-------|-----------|------------|----------|------------------|
| 1 | 100 | 150 | 0.05 | 95.29 \pm 0.04 |
| 50 | 100 | 150 | 0.05 | 95.10 \pm 0.02 |
| 100 | | | | 95.06 \pm 0.06 |
| 1 | 10 | 15 | 0.05 | 95.10 \pm 0.06 |
| | 300 | 450 | | 95.31 \pm 0.02 |
| 1 | 100 | 150 | 0.01 | 95.15 \pm 0.08 |
| | | | 0.10 | 95.24 \pm 0.03 |

6 CONCLUSION

533 In this work, we proposed A^2SG , a unified framework for training deep SNNs. By integrating
 534 spatio-temporal adaptation with a neuron-aware asymmetric design, A^2SG reduces gradient variabil-
 535 ity, stabilizes optimization, and encourages convergence to flatter minima. Our theoretical analysis
 536 establishes the link between gradient variation and loss landscape curvature. Moreover, we prove
 537 that the asymmetric surrogate achieves lower variation than its symmetric counterparts. Extensive
 538 experiments across diverse SNN architectures and tasks demonstrate that A^2SG consistently im-
 539 proves accuracy, robustness, and efficiency, highlighting surrogate gradient design as a key factor
 for reliable and scalable SNN training.

REFERENCES

- 540
541
542
543 Malyaban Bal and Abhronil Sengupta. Spikingbert: Distilling bert to train spiking language models
544 using implicit differentiation. In *Proceedings of the AAAI conference on artificial intelligence*,
545 volume 38, pp. 10998–11006, 2024.
546
- 547 Pratik Chaudhari, Anna Choromanska, Stefano Soatto, Yann LeCun, Carlo Baldassi, Christian
548 Borgs, Jennifer Chayes, Levent Sagun, and Riccardo Zecchina. Entropy-sgd: Biasing gradient
549 descent into wide valleys. *Journal of Statistical Mechanics: Theory and Experiment*, 2019(12):
550 124018, 2019.
- 551 Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated
552 data augmentation with a reduced search space. In *Proceedings of the IEEE/CVF conference on*
553 *computer vision and pattern recognition workshops*, pp. 702–703, 2020.
554
- 555 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hi-
556 erarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
557 pp. 248–255. Ieee, 2009.
- 558 Shikuang Deng, Yuhang Li, Shanghang Zhang, and Shi Gu. Temporal efficient training of spiking
559 neural network via gradient re-weighting. In *International Conference on Learning Representa-*
560 *tions*, 2022.
561
- 562 Wei Fang, Zhaofei Yu, Yanqi Chen, Tiejun Huang, Timothée Masquelier, and Yonghong Tian. Deep
563 residual learning in spiking neural networks. *Advances in Neural Information Processing Systems*,
564 34:21056–21069, 2021a.
- 565 Wei Fang, Zhaofei Yu, Yanqi Chen, Timothée Masquelier, Tiejun Huang, and Yonghong Tian. In-
566 corporating learnable membrane time constant to enhance learning of spiking neural networks.
567 In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 2661–2671,
568 2021b.
569
- 570 Behrooz Ghorbani, Shankar Krishnan, and Ying Xiao. An investigation into neural net optimization
571 via hessian eigenvalue density. In *Proceedings of the 36th International Conference on Machine*
572 *Learning (ICML)*, volume 97, pp. 2232–2241. PMLR, 2019.
- 573 Anne Greenbaum, Ren-cang Li, and Michael L Overton. First-order perturbation theory for eigen-
574 values and eigenvectors. *SIAM review*, 62(2):463–482, 2020.
575
- 576 Yufei Guo, Yuanpei Chen, Liwen Zhang, Xiaode Liu, Yinglei Wang, Xuhui Huang, and Zhe Ma. Im-
577 loss: information maximization loss for spiking neural networks. *Advances in Neural Information*
578 *Processing Systems*, 35:156–166, 2022.
- 579 Yufei Guo, Xiaode Liu, Yuanpei Chen, Liwen Zhang, Weihang Peng, Yuhan Zhang, Xuhui Huang,
580 and Zhe Ma. Rmp-loss: Regularizing membrane potential distribution for spiking neural net-
581 works. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp.
582 17391–17401, 2023a.
583
- 584 Yufei Guo, Yuhan Zhang, Yuanpei Chen, Weihang Peng, Xiaode Liu, Liwen Zhang, Xuhui Huang,
585 and Zhe Ma. Membrane potential batch normalization for spiking neural networks. In *Proceed-*
586 *ings of the IEEE/CVF International Conference on Computer Vision*, pp. 19420–19430, 2023b.
- 587 Yufei Guo, Yuanpei Chen, Zecheng Hao, Weihang Peng, Zhou Jie, Yuhan Zhang, Xiaode Liu, and
588 Zhe Ma. Take a shortcut back: Mitigating the gradient vanishing for training spiking neural
589 networks. *Advances in Neural Information Processing Systems*, 37:24849–24867, 2024.
590
- 591 Sepp Hochreiter and Jürgen Schmidhuber. Flat minima. *Neural computation*, 9(1):1–42, 1997.
592
- 593 Yangfan Hu, Huajin Tang, and Gang Pan. Spiking deep residual networks. *IEEE Transactions on*
Neural Networks and Learning Systems, 34(8):5200–5205, 2021.

- 594 Haiyan Jiang, Vincent Zoonekynd, Giulia De Masi, Bin Gu, and Huan Xiong. Tab: Temporal accu-
595 mulated batch normalization in spiking neural networks. In *The Twelfth International Conference*
596 *on Learning Representations*, 2024.
- 597
- 598 Jiaqiang Jiang, Lei Wang, Runhao Jiang, Jing Fan, and Rui Yan. Adaptive gradient learning for
599 spiking neural networks by exploiting membrane potential dynamics. In James Kwok (ed.), *Pro-*
600 *ceedings of the Thirty-Fourth International Joint Conference on Artificial Intelligence, IJCAI-25*,
601 pp. 4164–4172. International Joint Conferences on Artificial Intelligence Organization, 8 2025.
602 doi: 10.24963/ijcai.2025/464. Main Track.
- 603 R Karakida, S Akaho, and S-i Amari. Universal statistics of fisher information in deep neural
604 networks. *Neural Computation*, 2019.
- 605
- 606 Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Pe-
607 ter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. In
608 *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- 609 Dongheon Kim, Seungchan Oh, Jinwoo Lee, Sangwoo Seo, and Kwanghee Suh. Fisher-sam: Con-
610 necting sharpness-aware minimization and fisher information. *Proceedings of Machine Learning*
611 *Research*, 162:1234–1248, 2022a.
- 612
- 613 Seijoon Kim, Seongsik Park, Byunggook Na, Jongwan Kim, and Sungroh Yoon. Towards fast and
614 accurate object detection in bio-inspired spiking neural networks through bayesian optimization.
615 *IEEE Access*, 9:2633–2643, 2020a.
- 616 Seijoon Kim, Seongsik Park, Byunggook Na, and Sungroh Yoon. Spiking-yolo: spiking neural
617 network for energy-efficient object detection. *Proceedings of the AAAI conference on artificial*
618 *intelligence*, 34(07):11270–11277, 2020b.
- 619 Youngeun Kim, Joshua Chough, and Priyadarshini Panda. Beyond classification: Directly training
620 spiking neural networks for semantic segmentation. *Neuromorphic Computing and Engineering*,
621 2(4):044015, 2022b.
- 622
- 623 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
624 2009.
- 625
- 626 Zhenxin Lei, Man Yao, Jiakui Hu, Xinhao Luo, Yanye Lu, Bo Xu, and Guoqi Li. Spike2former:
627 Efficient spiking transformer for high-performance image segmentation. In *Proceedings of the*
628 *AAAI Conference on Artificial Intelligence*, volume 39, pp. 1364–1372, 2025.
- 629 Hongmin Li, Hanchao Liu, Xiangyang Ji, Guoqi Li, and Luping Shi. Cifar10-dvs: an event-stream
630 dataset for object classification. *Frontiers in neuroscience*, 11:244131, 2017.
- 631
- 632 Yuhang Li, Yufei Guo, Shanghang Zhang, Shikuang Deng, Yongqing Hai, and Shi Gu. Differ-
633 entiable spike: Rethinking gradient-descent for training spiking neural networks. *Advances in*
634 *Neural Information Processing Systems*, 34:23426–23439, 2021.
- 635
- 636 Shuang Lian, Jiangrong Shen, Qianhui Liu, Ziming Wang, Rui Yan, and Huajin Tang. Learnable
637 surrogate gradient for direct training spiking neural networks. In *Proceedings of the International*
Joint Conference on Artificial Intelligence, pp. 3002–3010, 2023.
- 638
- 639 Zhibin Liao, Tom Drummond, Ian Reid, and Gustavo Carneiro. Approximate fisher information
640 matrix to characterize the training of deep neural networks. *IEEE transactions on pattern analysis*
and machine intelligence, 42(1):15–26, 2018.
- 641
- 642 Hao Lin, Shikuang Deng, and Shi Gu. Efficient surrogate gradients for training spiking neural
643 networks. 2023.
- 644
- 645 W. Maass. Networks of spiking neurons: the third generation of neural network models. *Neural*
Networks, 10(9):1659–1671, 1997.
- 646
- 647 James Martens. New insights and perspectives on the natural gradient method. *Journal of Machine*
Learning Research, 21(146):1–76, 2020.

- 648 Emre O Neftci, Hesham Mostafa, and Friedemann Zenke. Surrogate gradient learning in spiking
649 neural networks: Bringing the power of gradient-based optimization to spiking neural networks.
650 *IEEE Signal Processing Magazine*, 36(6):51–63, 2019.
- 651
- 652 Seongsik Park, Seijoon Kim, Byungook Na, and Sungroh Yoon. T2fsnn: deep spiking neural
653 networks with time-to-first-spike coding. In *2020 57th ACM/IEEE design automation conference*
654 *(DAC)*, pp. 1–6. IEEE, 2020.
- 655 Qiaoyi Su, Yuhong Chou, Yifan Hu, Jianing Li, Shijie Mei, Ziyang Zhang, and Guoqi Li. Deep
656 directly-trained spiking neural networks for object detection. In *Proceedings of the IEEE/CVF*
657 *International Conference on Computer Vision*, pp. 6555–6565, 2023.
- 658
- 659 Amirhossein Tavanaei, Masoud Ghodrati, Saeed Reza Kheradpisheh, Timothée Masquelier, and
660 Anthony Maida. Deep learning in spiking neural networks. *Neural networks*, 111:47–63, 2019.
- 661 Yu-Lin Tsai, Chia-Yi Hsu, Chia-Mu Yu, and Pin-Yu Chen. Formalizing generalization and adver-
662 sarial robustness of neural networks to weight perturbations. *Advances in Neural Information*
663 *Processing Systems*, 34:19692–19704, 2021.
- 664 Yujie Wu, Lei Deng, Guoqi Li, and Luping Shi. Spatio-temporal backpropagation for training high-
665 performance spiking neural networks. *Frontiers in neuroscience*, 12:323875, 2018.
- 666
- 667 Xingrun Xing, Zheng Zhang, Ziyi Ni, Shitao Xiao, Yiming Ju, Siqi Fan, Yequan Wang, Jiajun
668 Zhang, and Guoqi Li. Spikelm: towards general spike-driven language modeling via elastic bi-
669 spiking mechanisms. In *Proceedings of the 41st International Conference on Machine Learning*,
670 *ICML’24*. JMLR.org, 2024.
- 671 Jiaqi Yan, Changping Wang, De Ma, Huajin Tang, Qian Zheng, and Gang Pan. Training high
672 performance spiking neural network by temporal model calibration. In *Forty-second International*
673 *Conference on Machine Learning*, 2025.
- 674
- 675 Man Yao, JiaKui Hu, Tianxiang Hu, Yifan Xu, Zhaokun Zhou, Yonghong Tian, Bo XU, and Guoqi
676 Li. Spike-driven transformer v2: Meta spiking neural network architecture inspiring the design
677 of next-generation neuromorphic chips. In *The Twelfth International Conference on Learning*
678 *Representations*, 2024.
- 679 Man Yao, Xuerui Qiu, Tianxiang Hu, Jiakui Hu, Yuhong Chou, Keyu Tian, Jianxing Liao, Luziwei
680 Leng, Bo Xu, and Guoqi Li. Scaling spike-driven transformer with efficient spike firing approxi-
681 mation training. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025.
- 682
- 683 Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo.
684 Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceed-*
685 *ings of the IEEE/CVF International Conference on Computer Vision*, pp. 6023–6032, 2019.
- 686 H Zhang, M Cisse, YN Dauphin, and D Lopez-Paz. mixup: beyondempirical risk minimization.
687 *iclr. Vancouver, BC, Canada, April*, 2018.
- 688
- 689 Dongcheng Zhao, Guobin Shen, Yiting Dong, Yang Li, and Yi Zeng. Improving stability and perfor-
690 mance of spiking neural networks through enhancing temporal consistency. *Pattern Recognition*,
691 159:111094, 2025.
- 692 Hanle Zheng, Yujie Wu, Lei Deng, Yifan Hu, and Guoqi Li. Going deeper with directly-trained
693 larger spiking neural networks. In *Proceedings of the AAAI conference on artificial intelligence*,
694 volume 35, pp. 11062–11070, 2021.
- 695
- 696 Xian Zhong, Shengwang Hu, Wenxuan Liu, Wenxin Huang, Jianhao Ding, Zhaofei Yu, and Tiejun
697 Huang. Towards low-latency event-based visual recognition with hybrid step-wise distillation
698 spiking neural networks. In *Proceedings of the 32nd ACM international conference on multime-*
699 *dia*, pp. 9828–9836, 2024.
- 700 Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene
701 parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and*
pattern recognition, pp. 633–641, 2017.

702 Zhaokun Zhou, Yuesheng Zhu, Chao He, Yaowei Wang, Shuicheng YAN, Yonghong Tian, and
703 Li Yuan. Spikformer: When spiking neural network meets transformer. In *International Confer-*
704 *ence on Learning Representations*, 2023.
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A APPENDIX

A.1 LEAKY INTEGRATE-AND-FIRE (LIF) NEURON

The dynamics of a LIF neuron can be formulated as follows. First, the membrane potential is updated at each time step according to:

$$u_i^l[t] = \frac{1}{\tau} \left(v_i^l[t-1] + \sum_j w_{ij}^l s_j^{l-1}[t] \right), \quad (\text{A1})$$

where u denotes the membrane potential, v the intermediate membrane state, w the synaptic weights, and s the input spike. Indices i and j represent the post- and pre-synaptic neurons, respectively, and l refers to the layer index. The parameters τ and t indicate the membrane time constant and discrete time step.

A spike is emitted when the membrane potential surpasses a predefined threshold:

$$s_i^l[t] = H(u_i^l[t] - V_{\text{th}}), \quad (\text{A2})$$

where $H(\cdot)$ is the Heaviside step function and V_{th} is the firing threshold.

After spike firing, the membrane potential is reset based on the intermediate state using the following mechanism:

$$v_i^l[t] = (u_i^l[t] - s_i^l[t]) s_i^l[t] + u_i^l[t] (1 - s_i^l[t]). \quad (\text{A3})$$

A.2 SPATIO-TEMPORAL BACKPROPAGATION (STBP)

SNNs require gradient propagation that considers not only spatial but also temporal variations. STBP is considered a suitable back-propagation method for SNNs as it incorporates both spatial and temporal components. The gradient of the loss L with respect to the membrane potential $u_i^l[t]$ is defined as follows:

$$\frac{\partial L}{\partial u_i^l[t]} = \frac{\partial L}{\partial s_i^l[t]} \frac{\partial s_i^l[t]}{\partial u_i^l[t]} + \frac{\partial L}{\partial u_i^l[t+1]} \frac{\partial u_i^l[t+1]}{\partial u_i^l[t]}, \quad (\text{A4})$$

here, $u[t]$ is the membrane potential at time t , and $s[t]$ is the spike counts at time t . The l and i indicate the layer and neuron index, respectively.

In addition, the gradient of the loss for the weight w^l in l th layer is given by:

$$\frac{\partial L}{\partial w^l} = \sum_t \frac{\partial L}{\partial w^l[t]} = \sum_t \frac{\partial L}{\partial u^l[t]} s^{l-1}[t]. \quad (\text{A5})$$

A.3 SURROGATE GRADIENTS IN STBP

Surrogate gradient replaces the non-differentiable function $\frac{\partial s}{\partial u}$ with smooth approximated functions to enable gradient-based training. Representative surrogate gradient functions include the boxcar function and the triangle function, which are defined as follows:

$$\frac{\partial s}{\partial u} = \frac{1}{2\beta} \cdot \mathbf{1}(|u - V_{\text{th}}| < \beta), \quad (\text{A6})$$

$$\frac{\partial s}{\partial u} = \frac{1}{\beta^2} \cdot \max(0, \beta - |u - V_{\text{th}}|). \quad (\text{A7})$$

The above equations represent the *BOX* function and the *TRI* function, respectively, and β is a parameter representing the effective window. These functions focus on mimicking the Dirac delta function, and their area remains constant regardless of the effective window. The function shapes are illustrated in Fig. A1

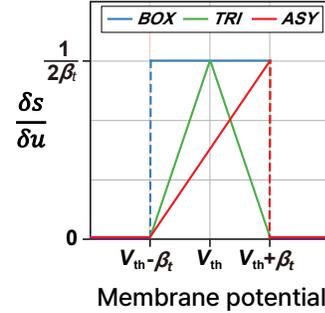


Figure A1: Visualization of surrogate gradient functions (*BOX*, *TRI*, *ASY*).

810 A.4 FISHER INFORMATION MATRIX (FIM)

811 The FIM is defined as:

$$812 \mathbf{F}(\theta) = \mathbb{E}_{(x,y) \sim p(x,y;\theta)} \left[\nabla_{\theta} \log p(y | x; \theta) \nabla_{\theta} \log p(y | x; \theta)^{\top} \right], \quad (\text{A8})$$

813 where θ denotes the model parameters, (x, y) is an input-output pair sampled from the data distri-
814 bution $p(x, y; \theta)$, and $p(y | x; \theta)$ is the model's conditional probability. The operator ∇_{θ} represents
815 the gradient with respect to θ . The FIM measures the variance of the gradient of the log-likelihood,
816 reflecting the sensitivity of the model parameters to changes in the data distribution. In this context,
817 the eigenvalues of the FIM quantify the curvature of the loss surface in various parameter directions.

821 A.5 PROOF OF THE DERIVATIVES OF SYMMETRIC SURROGATE GRADIENTS

822 Let $f : [\theta - \beta, \theta + \beta] \rightarrow \mathbb{R}_{\geq 0}$ be a symmetric surrogate function, with a unique maximum at $u = \theta$,
823 and

$$824 f(\theta - \beta) = f(\theta + \beta) = m (\geq 0), \quad \int_{\theta - \beta}^{\theta + \beta} f(u) du = c > 0.$$

825 Set the *effective area above the floor* by

$$826 c_{\text{eff}} := \int_{\theta - \beta}^{\theta + \beta} (f(u) - m) du = c - 2\beta m.$$

827 We assume $c_{\text{eff}} > 0$ (otherwise $f \equiv m$ on the window and all derivatives vanish).

828 Define $g(u) := f(u) - m$. Then g is symmetric, $g(\theta \pm \beta) = 0$, and $\int_{\theta - \beta}^{\theta + \beta} g(u) du = c_{\text{eff}}$. Let
829 $M_f := f(\theta)$ and $M_g := g(\theta) = M_f - m$.

830 **Height.** By symmetry,

$$831 c_{\text{eff}} = 2 \int_{\theta}^{\theta + \beta} g(u) du \leq 2\beta M_g \Rightarrow M_g \geq \frac{c_{\text{eff}}}{2\beta},$$

832 hence

$$833 M_f = m + M_g \geq m + \frac{c_{\text{eff}}}{2\beta} = \frac{c}{2\beta}.$$

834 In particular, $\|f\|_{\infty} \geq M_f \geq c/(2\beta)$.

835 **First derivative.** On $[\theta, \theta + \beta]$, $g(\theta) = M_g$ and $g(\theta + \beta) = 0$. By the mean value theorem there
836 exists ξ with

$$837 |g'(\xi)| = \frac{M_g}{\beta}.$$

838 Since $g' = f'$, we obtain

$$839 \|f'\|_{\infty} \geq \frac{M_g}{\beta} \geq \frac{c_{\text{eff}}}{2\beta^2} = \frac{c - 2\beta m}{2\beta^2}.$$

840 **Conclusion.** With $c_{\text{eff}} = c - 2\beta m > 0$,

$$841 \|f\|_{\infty} \geq \frac{c}{2\beta} = \Omega(\beta^{-1}), \quad \|f'\|_{\infty} \geq \frac{c - 2\beta m}{2\beta^2} = \Omega(\beta^{-2}).$$

842 Substituting into

$$843 \frac{\partial^2 L}{\partial w^2} = \frac{\partial^2 L}{\partial H^2} (H'(u)x)^2 + \frac{\partial L}{\partial H} H''(u) x^2 = \frac{\partial^2 L}{\partial H^2} (f(u)x)^2 + \frac{\partial L}{\partial H} f'(u) x^2, \quad (\text{A9})$$

844 we obtain

$$845 \left| \frac{\partial^2 L}{\partial w^2} \right|_{SNN} = \Omega\left(\frac{x^2}{\beta^2}\right).$$

846 which shows sharper curvature than the $\mathcal{O}(x^2)$ scaling of smooth DNNs.

A.6 TRI INDUCES LARGER CURVATURE THAN BOX UNDER AREA NORMALIZATION

Setup. Consider the surrogate gradients supported on $[\theta - \beta, \theta + \beta]$ with unit area. The boxcar (*BOX*) and triangular (*TRI*) surrogates are given by Eqs. A6 and A7. For a weight w with input x and pre-activation $u = wx$, the Hessian contribution is stated in Eq. A9.

BOX properties. The *BOX* function has support length 2β , constant height $1/(2\beta)$, and area 1. Therefore,

$$\|f_{\text{BOX}}\|_{\infty} = \frac{1}{2\beta}, \quad \|f'_{\text{BOX}}\| = 0 \text{ almost everywhere on } (\theta - \beta, \theta + \beta).$$

TRI properties. The *TRI* function is piecewise linear with the maximum at $u = \theta$ given by $f_{\text{TRI}}(\theta) = 1/\beta$, twice the peak of *BOX* under the same area constraint. The slopes are $\pm 1/\beta^2$, hence

$$\|f_{\text{TRI}}\|_{\infty} = \frac{1}{\beta}, \quad \|f'_{\text{TRI}}\|_{\infty} = \frac{1}{\beta^2}.$$

Comparison. Relative to *BOX*,

$$\|f_{\text{TRI}}\|_{\infty} = 2 \|f_{\text{BOX}}\|_{\infty}, \quad \|f'_{\text{TRI}}\|_{\infty} > \|f'_{\text{BOX}}\|_{\infty}.$$

Thus, *TRI* attains both higher peak and steeper slope.

Implication for curvature. Substituting into Eq. A9, the first term scales with $f(u)^2$ and is therefore at least four times larger for *TRI* near $u = \theta$ compared to *BOX*. The second term involves $f'(u)$, which is strictly larger for *TRI* inside the window since $f'_{\text{BOX}} = 0$ almost everywhere. As a result, both contributions to the curvature are enhanced under *TRI*, leading to the following conclusion:

$$\left| \frac{\partial^2 L}{\partial w^2} \right|_{\text{TRI}} \geq 2 \cdot \left| \frac{\partial^2 L}{\partial w^2} \right|_{\text{BOX}}$$

A.7 ADAPTIVE SURROGATE GRADIENTS CONSIDERING SGV AND TGC

Algorithm A1 Adaptive Surrogate Gradients Considering SGV and TGC

```

1: Input: iteration  $i$ ,  $\frac{\partial L}{\partial s[t]}$ 
2: Output:  $\frac{\partial L}{\partial u[t]}$ 
3: for  $l = L$  to 1 do
4:   if  $i = 0$  then
5:      $\beta_{[t:0,\dots,T]}^l \leftarrow \beta_{\text{init}}$ 
6:   end if
7:   if  $((i \bmod i_{\text{update}}) = 0)$  and  $(i \neq 0)$  then
8:     for  $t = T$  to 1 do
9:       if  $t = T$  then
10:         $\beta_t^l \leftarrow \beta_{\text{search}}(\text{SGV}, \frac{\partial L}{\partial s[t]}, u^l[t])$ 
11:       else
12:         $\beta_t^l \leftarrow \beta_{\text{search}}(\text{TGC}, g_{\text{cur}}^l[t], \frac{\partial L}{\partial s[t]}, u^l[t])$ 
13:       end if
14:     end for
15:   else
16:     for  $t = T$  to 1 do
17:        $g_{\text{cur}}^l[t] \leftarrow f(u^l[t], \beta_t^l)$ 
18:     end for
19:   end if
20: end for

```

Alg. A1 outlines the adaptive surrogate gradient procedure that simultaneously considers both SGV and TGC. At the beginning of training, the effective window β is initialized. During every update step, β is adaptively adjusted according to the observed gradient distribution: SGV is used to calibrate the final timestep, while TGC is employed for earlier timesteps to maintain temporal consistency of error signals. When not in an update step, the algorithm computes the current layer gradient using the stored β . In this way, the method dynamically tunes β across layers and timesteps.

A.8 BETA SEARCH VIA BAYESIAN OPTIMIZATION

Algorithm A2 Beta Search via Bayesian Optimization

```

922 1: Input: metric  $M$ ,  $g_{\text{ref}}$ ,  $\frac{\partial L}{\partial s[t]}$ ,  $u[t]$ ,  $t$ ,  $j$ 
923 2: Output:  $\beta$ 
924 3:  $\beta_{\text{min}} \leftarrow \beta_{\text{min,init}}$ ,  $\beta_{\text{max}} \leftarrow \beta_{\text{max,init}}$ 
925 4:  $\beta_{\text{obs}} \leftarrow \text{Random}(\beta_{\text{min}}, \beta_{\text{max}}, n_{\text{obs}})$ 
926 5: if  $t = T$  then
927 6:    $M \leftarrow \text{SGV}(\beta_{\text{obs}}, u[t], \frac{\partial L}{\partial s[t]})$ 
928 7: else
929 8:    $M \leftarrow \text{TGC}(\beta_{\text{obs}}, g_{\text{ref}}, u[t], \frac{\partial L}{\partial s[t]})$ 
930 9: end if
931 10:  $f_{\text{best}} \leftarrow \max(M)$ 
932 11:  $\beta_{\text{best}} \leftarrow \beta_{\text{obs}}[\text{argmax}(M)]$ 
933 12:  $\beta_{\text{min}} \leftarrow \max(\beta_{\text{best}} - \delta, \beta_{\text{min}})$ 
934 13:  $\beta_{\text{max}} \leftarrow \min(\beta_{\text{best}} + \delta, \beta_{\text{max}})$ 
935 14:  $\beta_{\text{eval}} \leftarrow \text{Uniform}(\beta_{\text{min}}, \beta_{\text{max}}, n_{\text{eval}})$ 
936 15:  $\mu_s, \sigma_s \leftarrow \text{GP}(\beta_{\text{obs}}, M, \beta_{\text{eval}})$ 
937 16:  $\text{EI} \leftarrow \text{ExpectedImprovement}(\mu_s, \sigma_s, f_{\text{best}})$ 
938 17:  $\beta^* \leftarrow \text{arg max}_{\beta} \text{EI}$ 
939 18: if  $\beta^* < f_{\text{best}}$  then
940 19:    $\beta^* \leftarrow \beta_{\text{best}}$ 
941 20: end if
942 21:  $\beta \leftarrow \beta^*$ 

```

Alg. A2 describes the adaptive search strategy for the effective window β . The method initializes a search range and samples candidates to evaluate the training metric. Depending on the current timestep t , either SGV or TGC is used to compute the evaluation metric. A Gaussian Process (GP) surrogate model is then fitted over the observed results, and the expected improvement (EI) criterion guides the selection of the next candidate β . The algorithm iteratively narrows down the search interval around the best candidate, ensuring efficient exploration while avoiding unstable regions. The final β is set to the candidate with the highest improvement score.

A.9 COMPUTING INFRASTRUCTURE

All experiments are performed on servers equipped with Intel(R) Xeon(R) Gold 6226R CPUs (2.90GHz, 520GB RAM) and NVIDIA RTX A6000 GPUs (8 units), running Ubuntu 20.04. Our implementation is based on CUDA 11.7, PyTorch 2.0.1 for ImageNet, Tensorflow/Keras 2.11.0 for CIFAR-10, CIFAR-100, and CIFAR10-DVS.

A.10 EXPERIMENTAL SETUP

The input size of the model is set to 32x32 for CIFAR10/100, 224x224 for ImageNet, and 48x48 for CIFAR10-DVS. On CIFAR10/100, We trained each model for 310 epochs with the AdamW optimizer and a cosine decay learning rate scheduler with a 20-epoch warm-up. We set it to 200 epochs for CIFAR10-DVS. All models on CIFAR10 and CIFAR10-DVS were trained with an initial learning rate of 1×10^{-5} , a learning rate of 6×10^{-3} , and a weight decay of 2×10^{-2} . On CIFAR100, all models were trained with an initial learning rate of 1×10^{-4} , a learning rate of 5×10^{-3} , and a weight decay of 4×10^{-2} . Data augmentation was performed using a combination of CutMix Yun et al. (2019) and RandAugment Cubuk et al. (2020). RandAugment was configured with one augmentation per image, a magnitude of 1, a magnitude standard deviation of 0.4, and an application rate of 0.5. The batch size was set to 100 for CIFAR10/100. For ImageNet experiments, we employed the E-Spikeformer (Yao et al., 2025) model. A batch size of 360 was used for the 10M model, while a batch size of 100 was used for the 173M model. Both models were trained with a base learning rate of 6×10^{-4} , a minimum learning rate of 1×10^{-6} , and 5 warm-up epochs. For data augmentation, we applied Mixup (Zhang et al., 2018) with a weight of 0.8 and CutMix with a

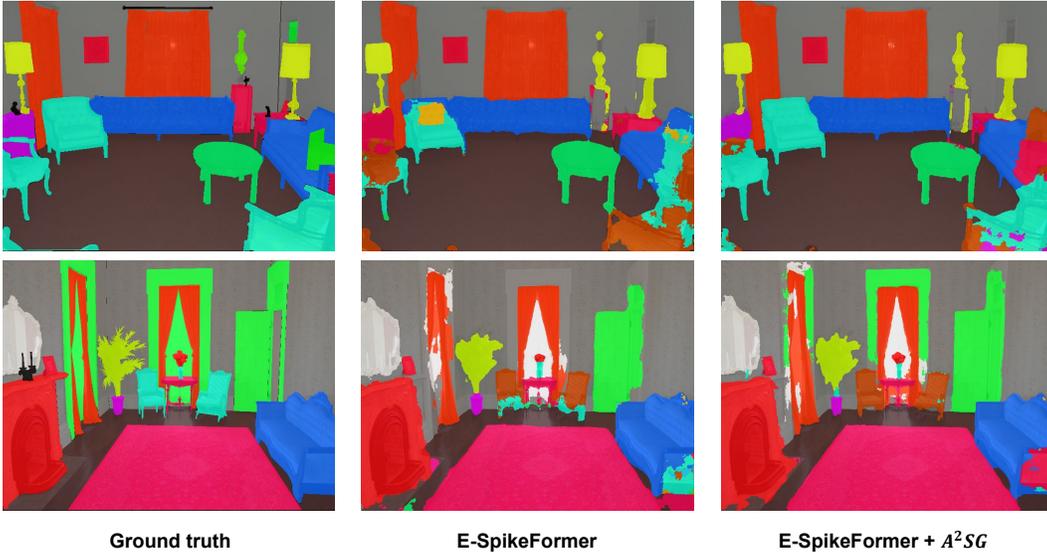


Figure A2: Segmentation results on the ADE20K dataset.

weight of 1.0. For all experiments, the effective window β was initialized to 0.5, and the optimal β during training was searched at the first iteration of every epoch.

A.11 VISUALIZATION OF SEGMENTATION

Fig. A2 presents segmentation results on the ADE20K dataset. Compared to the E-spikeFormer, the proposed A^2SG improves boundary sharpness and object consistency.

A.12 QUANTITATIVE ANALYSIS OF THE ASYMMETRIC SURROGATE GRADIENT

We provide additional analysis of the effect of ASY on VGG16 using CIFAR-10. As shown in Table A1, in the Conv1 layer, ASY substantially increases the proportion of silent neurons and reduces the total spike count (7.92k), outperforming both BOX (9.12k) and TRI (9.76k), while also yielding higher test accuracy.

To further interpret this result, we analyze the membrane potential statistics.

In SNNs, the weight distribution is mapped linearly to each neuron’s membrane potential. Given an input spike count M , the membrane potential u can be expressed as

$$u = \sum_i w_i x_i,$$

where x_i and w_i denote the presynaptic spikes and their corresponding weights, respectively. The mean and standard deviation of the membrane potential are $\mu_u = M\mu_w$ and $\sigma_u = \sqrt{M}\sigma_w$, where μ_w and σ_w are the mean and standard deviation of the weight distribution. By the central limit theorem, the membrane potential can be approximated by $\mathcal{N}(\mu_u, \sigma_u^2)$. To quantify the distance between the mean membrane potential and the firing threshold V_{th} , we define the relative membrane distance (RMD) as

$$\text{RMD} := \frac{V_{th} - \mu_u}{\sigma_u}. \quad (\text{A10})$$

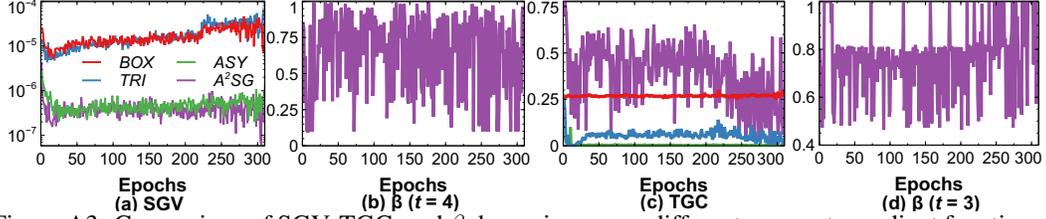
A larger RMD indicates that the membrane potential lies farther below the threshold, thus reducing the probability of spiking.

Table A1: Distribution of spikes in the Conv1 layer for BOX , TRI , and ASY . Each value represents the percentage of samples with the corresponding spike counts.

| Spike Counts | BOX | TRI | ASY |
|----------------------|--------------|--------------|--------------|
| 0 | 91.73% | 92.41% | 93.24% |
| 1 | 6.35% | 5.58% | 5.15% |
| 2 | 1.49% | 1.53% | 1.25% |
| 3 | 0.29% | 0.32% | 0.24% |
| 4 | 0.14% | 0.16% | 0.12% |
| Num of Spikes | 9.12k | 9.76k | 7.92k |

Table A2: μ_w , σ_w , and RMD depending on the surrogate gradient (SG) functions in Conv1.

| SG func. | $\mu_w (\times 10^{-3})$ | σ_w | RMD | # of spikes ($\times 10^3$) |
|----------|--------------------------|------------|-------|-------------------------------|
| BOX | -8.122 | 0.1408 | 7.159 | 9.12 |
| TRI | -8.205 | 0.1394 | 7.231 | 9.76 |
| ASY | -7.002 | 0.1361 | 7.398 | 7.92 |

Figure A3: Comparison of SGV, TGC, and β dynamics across different surrogate gradient functions at Conv1 layer. (a) SGV over epochs, (b) β dynamics at $t = 4$, (c) TGC over epochs, (d) β dynamics at $t = 3$.

As summarized in Table A2, *ASY* achieves the highest RMD and the lowest spike count, confirming its effectiveness in suppressing redundant spikes and improving energy efficiency.

A.13 COMPARISON OF SGV, TGC, AND β DYNAMICS ON OTHER LAYERS

Fig. A3 and A4 illustrate SGV, TGC, and β for Conv1 and Conv3, respectively. Similar to the results in Conv5-2, *A²SG* achieves the lowest SGV and the highest TGC. Notably, while *ASY* exhibited relatively high TGC in Conv5-2, it shows much lower TGC in the earlier layers. In contrast, *A²SG* consistently maintains high TGC across all layers by applying the adaptive surrogate strategy to *ASY*. Furthermore, in Figs. A3 and A4-(b), (d), it can be observed that β is adjusted in a manner that improves both SGV and TGC.

A.14 THEORETICAL ANALYSIS

Theorem A1 (CV-Minimizing Symmetric Function under Area and Boundary Constraints). *Let $f : [a, b] \rightarrow \mathbb{R}_{\geq 0}$ be a function satisfying, where $a = \theta - \beta$ and $b = \theta + \beta$:*

- (**Symmetry**): $f(u) = f(a + b - u)$ for all $u \in [a, b]$
- (**Boundary condition**): $f(a) = f(b) \geq 0$
- (**Nonnegativity**): $f(u) \geq 0$
- (**Area constraint**): $\int_a^b f(u) du = c > 0$

Let the weight function be given by the Gaussian density

$$p(u) := \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(u-\mu)^2}{2\sigma^2}\right), \quad \text{with } \mu < a,$$

so that $p(u)$ is strictly decreasing and convex on $[a, b]$.

Then the unique function f^* that minimizes the CV

$$\text{CV}_{\text{sim}}[f] := \frac{\sqrt{\int_a^b f(u)^2 p(u) du - \left(\int_a^b f(u) p(u) du\right)^2}}{\int_a^b f(u) p(u) du},$$

over all such admissible functions is the symmetric triangular function

$$f^*(u) := \frac{4c}{(b-a)^2} \cdot \min(u-a, b-u).$$

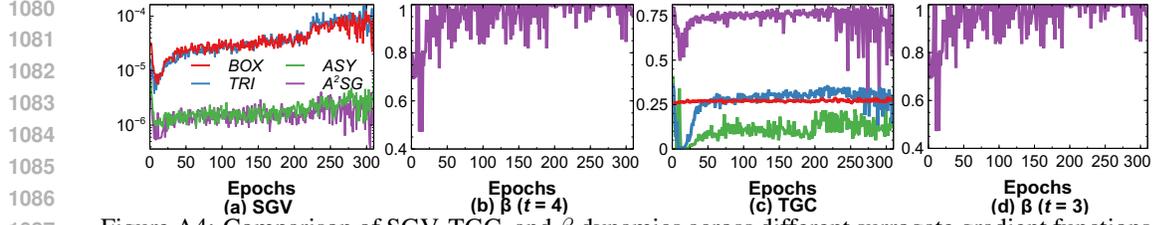


Figure A4: Comparison of SGV, TGC, and β dynamics across different surrogate gradient functions at Conv3 layer. (a) SGV over epochs, (b) β dynamics at $t = 4$, (c) TGC over epochs, (d) β dynamics at $t = 3$.

Proof. Let $m = \frac{a+b}{2}$ be the midpoint of the interval. Due to the symmetry condition, any admissible function f is uniquely determined by its restriction g to $[a, m]$, with the reconstruction:

$$f(u) = \begin{cases} g(u), & u \in [a, m], \\ g(a+b-u), & u \in [m, b]. \end{cases}$$

Since $p(u)$ is strictly decreasing and convex on $[a, b]$, it is also decreasing on $[a, m]$. Over this domain, define the following Rayleigh-type quotient:

$$R[g] := \frac{\int_a^m g(u)^2 p(u) du + \int_m^b g(a+b-u)^2 p(u) du}{\left(\int_a^m g(u) p(u) du + \int_m^b g(a+b-u) p(u) du \right)^2}.$$

By change of variable $u' = a+b-u$, and using symmetry of f , we have:

$$R[g] = \frac{2 \int_a^m g(u)^2 p(u) du}{\left(2 \int_a^m g(u) p(u) du \right)^2} = \frac{\int_a^m g(u)^2 p(u) du}{\left(\int_a^m g(u) p(u) du \right)^2}.$$

This Rayleigh quotient is minimized when $g(u)$ is proportional to a linear function increasing from a to m . That is:

$$g^*(u) = \alpha(u-a), \quad u \in [a, m],$$

with boundary condition $g(a) = 0$. Extending symmetrically gives:

$$f^*(u) = \alpha \cdot \min(u-a, b-u).$$

To satisfy the area constraint:

$$\int_a^b f^*(u) du = \int_a^m \alpha(u-a) du + \int_m^b \alpha(b-u) du = \alpha \cdot \frac{(b-a)^2}{4} = c,$$

which yields:

$$\alpha = \frac{4c}{(b-a)^2}.$$

Thus, the minimizing function is:

$$f^*(u) = \frac{4c}{(b-a)^2} \cdot \min(u-a, b-u),$$

which is the symmetric triangular function. \square

Theorem A2 (CV Comparison of Asymmetric and Symmetric Surrogates). *Let $f_{\text{asy}}(u)$ and $f_{\text{sym}}(u)$ be asymmetric and symmetric surrogate gradient functions defined over $[a, b]$, satisfying the boundary condition $f(a) = f(b) = 0$, nonnegativity $f(u) \geq 0$, and area constraint $\int_a^b f(u) du = c$, where $a = \theta - \beta$ and $b = \theta + \beta$. Suppose the membrane potential $u \sim \mathcal{N}(\mu, \sigma^2)$ with $\mu < a$, so that $p(u)$ is strictly decreasing on $[a, b]$. Then, under a linear approximation of the Gaussian, we have:*

$$\text{CV}_{\text{asy}} < \text{CV}_{\text{sym}} \quad \text{if } L\kappa > \sigma^2.$$

1134 *Proof.* We define $L := b - a$ and approximate the Gaussian by a first-order Taylor expansion around
 1135 $u = a$:

$$1136 \quad A := p(a) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(a - \mu)^2}{2\sigma^2}\right),$$

1137

$$1138 \quad B := -p'(a) = \frac{a - \mu}{\sigma^2} \cdot A.$$

1139 Let $\kappa := a - \mu$, so that:

$$1140 \quad A = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{\kappa^2}{2\sigma^2}}, \quad B = \frac{\kappa}{\sigma^2} A.$$

1141 Under the area constraint, define the surrogates:

$$1142 \quad f_{\text{asy}}(u) = \frac{2c}{L} \cdot \frac{u - a}{L},$$

$$1143 \quad \text{and } f_{\text{sym.tri}}(u) = \begin{cases} \frac{2c}{L} \cdot \frac{u-a}{L/2}, & u \leq \theta, \\ \frac{2c}{L} \cdot \frac{b-u}{L/2}, & u > \theta, \end{cases} \quad \text{where } \theta = \frac{a+b}{2}.$$

1144 **Asymmetric surrogate:**

$$1145 \quad f_{\text{asy}}(u) = \frac{2c}{L^2}(u - a)$$

1146 *Expectation:*

$$1147 \quad \mathbb{E}[f_{\text{asy}}] = \int_a^b f_{\text{asy}}(u) \cdot p(u) du$$

$$1148 \quad = \frac{2c}{L^2} \int_a^b (u - a)(A - B(u - a)) du$$

$$1149 \quad = \frac{2c}{L^2} \int_0^L x(A - Bx) dx$$

$$1150 \quad = \frac{2c}{L^2} \left(A \cdot \frac{L^2}{2} - B \cdot \frac{L^3}{3} \right)$$

$$1151 \quad = c \left(A - \frac{2}{3}BL \right)$$

1152 *Second moment:*

$$1153 \quad \mathbb{E}[f_{\text{asy}}^2] = \left(\frac{2c}{L^2} \right)^2 \int_0^L x^2(A - Bx) dx$$

$$1154 \quad = \frac{4c^2}{L^4} \left(A \cdot \frac{L^3}{3} - B \cdot \frac{L^4}{4} \right)$$

$$1155 \quad = \frac{4c^2}{L} \left(\frac{A}{3} - \frac{BL}{4} \right)$$

1156 *Variance:*

$$1157 \quad \text{Var}[f_{\text{asy}}] = \mathbb{E}[f_{\text{asy}}^2] - \mathbb{E}[f_{\text{asy}}]^2$$

$$1158 \quad = \frac{4c^2}{L} \left(\frac{A}{3} - \frac{BL}{4} \right) - c^2 \left(A - \frac{2}{3}BL \right)^2$$

1159 **Symmetric surrogate (triangle function):**

$$1160 \quad f_{\text{sym.tri}}(u) = \begin{cases} \frac{4c}{L^2}(u - a), & u \in [a, \theta], \\ \frac{4c}{L^2}(b - u), & u \in [\theta, b] \end{cases}, \quad \theta = \frac{a+b}{2}$$

1188 *Expectation:*

$$\begin{aligned}
 1189 \mathbb{E}[f_{\text{sym.tri}}] &= \frac{4c}{L^2} \left[\int_a^\theta (u-a)p(u) du + \int_\theta^b (b-u)p(u) du \right] \\
 1190 &= \frac{8c}{L^2} \int_0^{L/2} x(A-Bx) dx \\
 1191 &= \frac{8c}{L^2} \left(A \cdot \frac{(L/2)^2}{2} - B \cdot \frac{(L/2)^3}{3} \right) \\
 1192 &= c \left(A - \frac{1}{3}BL \right)
 \end{aligned}$$

1200 *Second moment:*

$$\begin{aligned}
 1201 \mathbb{E}[f_{\text{sym.tri}}^2] &= \frac{8c^2}{L^4} \int_0^{L/2} x^2(A-Bx) dx \cdot 2 \\
 1202 &= \frac{16c^2}{L^4} \left(A \cdot \frac{(L/2)^3}{3} - B \cdot \frac{(L/2)^4}{4} \right) \\
 1203 &= \frac{2c^2}{L} \left(\frac{A}{3} - \frac{BL}{8} \right)
 \end{aligned}$$

1209 *Variance:*

$$\text{Var}[f_{\text{sym.tri}}] = \frac{2c^2}{L} \left(\frac{A}{3} - \frac{BL}{8} \right) - c^2 \left(A - \frac{1}{3}BL \right)^2$$

1213 **Coefficient Ratio:**

$$\begin{aligned}
 1214 r &:= \frac{\mathbb{E}[f_{\text{asy}}]}{\mathbb{E}[f_{\text{sym.tri}}]} = \frac{A - \frac{2}{3}BL}{A - \frac{1}{3}BL} = \frac{3A - 2BL}{3A - BL} \\
 1215 & \\
 1216 & \\
 1217 \eta &:= \frac{\text{Var}[f_{\text{asy}}]}{\text{Var}[f_{\text{sym.tri}}]} = \frac{\frac{4c^2}{L} \left(\frac{A}{3} - \frac{BL}{4} \right) - c^2 \left(A - \frac{2}{3}BL \right)^2}{\frac{2c^2}{L} \left(\frac{A}{3} - \frac{BL}{8} \right) - c^2 \left(A - \frac{1}{3}BL \right)^2}
 \end{aligned}$$

1220 Then the squared CV ratio becomes:

$$\begin{aligned}
 1221 \left(\frac{\text{CV}_{\text{asy}}}{\text{CV}_{\text{sym.tri}}} \right)^2 &= \frac{\eta}{r^2} \\
 1222 &= \frac{3(2A - BL)^2 \cdot [12A - 9BL + L(3A - 2BL)^2]}{(3A - 2BL)^2 \cdot [16A - 8BL + 3L(2A - BL)^2]},
 \end{aligned}$$

1227 which simplifies under substitution to:

$$\frac{\eta}{r^2} \sim \frac{3}{4} \cdot \left(\frac{L\kappa - 2\sigma^2}{2L\kappa - 3\sigma^2} \right)^2$$

1231 and thus yields the condition $\text{CV}_{\text{asy}} < \text{CV}_{\text{sym.tri}}$ if $L\kappa > \sigma^2$. Therefore, according to the Theo-
 1232 rem 1, $\text{CV}_{\text{asy}} < \text{CV}_{\text{sym}}$ if $L\kappa > \sigma^2$.

1233 \square

1235 A.15 EXTENSION OF THEOREM 2 WITH n -SEGMENT PIECEWISE-LINEAR APPROXIMATION

1237 In this subsection, we show that the conclusion of Theorem 2 is stable when the Gaussian weight
 1238 function on the effective window is approximated by an n -segment piecewise-linear model instead of
 1239 a single linear segment. We use the same notation as in Theorem 2 and Theorem A2: the effective
 1240 window is $I = [a, b] = [\theta - \beta, \theta + \beta]$ with width $L = b - a$, the membrane potential satisfies
 1241 $u \sim \mathcal{N}(\mu, \sigma^2)$ with $\mu < a$, and $p(u)$ denotes the Gaussian weight function. We also define $\delta :=$
 $a - \mu > 0$.

One-segment and n -segment PWLA models. We now specify the weight-function models on I .

- **One-segment linear model \tilde{p}_1 .** Following Theorem 2 (and Theorem A2), we approximate $p(u)$ on I by its first-order Taylor expansion at the left boundary $u = a$:

$$\tilde{p}_1(u) := p(a) + p'(a)(u - a), \quad u \in I.$$

- **n -segment piecewise-linear model \tilde{p}_n .** For $n \in \mathbb{N}$, we divide I into n equal sub-intervals of length $\Delta := L/n$ with grid points $u_k := a + k\Delta$ for $k = 0, \dots, n$. The n -segment piecewise-linear approximation of p is defined by

$$\tilde{p}_n(u) := p(u_k) + p'(u_k)(u - u_k) \quad \text{for } u \in [u_k, u_{k+1}), \quad k = 0, \dots, n-1. \quad (\text{A11})$$

We denote $R_1^2 := R(\tilde{p}_1)^2$ and $R_n^2 := R(\tilde{p}_n)^2$.

We denote by \mathcal{H} the class of weight functions considered in this paper, namely the Gaussian weight $p(u)$ and its PWLA approximations \tilde{p}_1, \tilde{p}_n on the effective window $I = [\theta - \beta, \theta + \beta]$. We refer to elements of \mathcal{H} as admissible weight functions.

Lemma 1 (Bound on the variation of $R(h)^2$). *Let $I = [a, b]$ be the effective window. For any nonnegative weight function $h : I \rightarrow \mathbb{R}_{\geq 0}$, define*

$$\begin{aligned} M_1(h) &= \int_a^b f_{\text{asy}}(u) h(u) du, & M_2(h) &= \int_a^b f_{\text{asy}}(u)^2 h(u) du, \\ J_1(h) &= \int_a^b f_{\text{sym}}(u) h(u) du, & J_2(h) &= \int_a^b f_{\text{sym}}(u)^2 h(u) du. \end{aligned}$$

Based on these, we define

$$\begin{aligned} \text{CV}_{\text{asy}}(h)^2 &:= \frac{M_2(h) - M_1(h)^2}{M_1(h)^2}, \\ \text{CV}_{\text{sym}}(h)^2 &:= \frac{J_2(h) - J_1(h)^2}{J_1(h)^2}, \\ R(h)^2 &:= \left(\frac{\text{CV}_{\text{asy}}(h)}{\text{CV}_{\text{sym}}(h)} \right)^2. \end{aligned}$$

Assume that there exist constants $m_{\min} > 0$ and $v_{\min} > 0$ such that, for all admissible weight functions h in the class considered in this paper,

$$M_1(h), J_1(h) \geq m_{\min}, \quad M_2(h) - M_1(h)^2, J_2(h) - J_1(h)^2 \geq v_{\min}. \quad (\text{A12})$$

Then there exists a constant $K > 0$ such that, for any admissible h, \tilde{h} ,

$$|R(h)^2 - R(\tilde{h})^2| \leq K \|h - \tilde{h}\|_{\infty}, \quad (\text{A13})$$

where $\|h - \tilde{h}\|_{\infty} := \sup_{u \in I} |h(u) - \tilde{h}(u)|$.

Proof. We denote $R(h)^2$ as

$$R(h)^2 = \Phi(v(h)),$$

where

$$v(h) := (M_1(h), M_2(h), J_1(h), J_2(h)) \in \mathbb{R}^4,$$

and

$$\Phi(m_1, m_2, j_1, j_2) := \frac{j_1^2}{m_1^2} \cdot \frac{m_2 - m_1^2}{j_2 - j_1^2}.$$

Step 1: Sensitivity of $v(h)$ with respect to h . Let $L := b - a$ and let h, \tilde{h} be two admissible weight functions. Set $\Delta h(u) := h(u) - \tilde{h}(u)$. Since f_{asy} and f_{sym} are fixed surrogate functions on the effective window I , they are bounded:

$$C_1 := \sup_{u \in I} (|f_{\text{asy}}(u)|, |f_{\text{sym}}(u)|) < \infty, \quad C_2 := \sup_{u \in I} (f_{\text{asy}}(u)^2, f_{\text{sym}}(u)^2) < \infty.$$

1296 We first bound the change in M_1 :

$$1297$$

$$1298 \quad M_1(h) - M_1(\tilde{h}) = \int_a^b f_{\text{asy}}(u) h(u) du - \int_a^b f_{\text{asy}}(u) \tilde{h}(u) du$$

$$1299$$

$$1300 \quad = \int_a^b f_{\text{asy}}(u) \Delta h(u) du.$$

$$1301$$

$$1302$$

1303 Using the triangle inequality and the definition of the sup norm,

$$1304$$

$$1305 \quad |M_1(h) - M_1(\tilde{h})| \leq \int_a^b |f_{\text{asy}}(u)| |\Delta h(u)| du$$

$$1306$$

$$1307 \quad \leq \left(\sup_{u \in I} |f_{\text{asy}}(u)| \right) \int_a^b |\Delta h(u)| du$$

$$1308$$

$$1309 \quad \leq C_1 \cdot L \cdot \sup_{u \in I} |\Delta h(u)|$$

$$1310$$

$$1311 \quad = C_1 L \|h - \tilde{h}\|_\infty.$$

$$1312$$

1313 Similarly, for M_2 we have

$$1314$$

$$1315 \quad M_2(h) - M_2(\tilde{h}) = \int_a^b f_{\text{asy}}(u)^2 \Delta h(u) du,$$

$$1316$$

1317 and hence

$$1318 \quad |M_2(h) - M_2(\tilde{h})| \leq C_2 L \|h - \tilde{h}\|_\infty.$$

1319 The same argument with f_{sym} and f_{sym}^2 yields

$$1320$$

$$1321 \quad |J_1(h) - J_1(\tilde{h})| \leq C_1 L \|h - \tilde{h}\|_\infty, \quad |J_2(h) - J_2(\tilde{h})| \leq C_2 L \|h - \tilde{h}\|_\infty.$$

$$1322$$

1323 Therefore, there exists a constant $C_v > 0$ ($C_v := L \max\{C_1, C_2\}$) such that

$$1324$$

$$1325 \quad \|v(h) - v(\tilde{h})\| \leq C_v \|h - \tilde{h}\|_\infty, \tag{A14}$$

$$1326$$

1327 where $\|\cdot\|$ denotes the Euclidean norm on \mathbb{R}^4 .

1328 **Step 2: Sensitivity of $\Phi(v)$ with respect to v .** By the assumptions in equation A12, for admissible

1329 weight functions h , we have

$$1330 \quad M_1(h), J_1(h) \geq m_{\min}, \quad M_2(h) - M_1(h)^2, J_2(h) - J_1(h)^2 \geq v_{\min}.$$

$$1331$$

$$1332$$

1333 Moreover, since the surrogates and weight functions are bounded on the finite interval I , the integrals

1334 $M_1(h), M_2(h), J_1(h), J_2(h)$ are bounded. Hence all vectors $v(h)$ lie in a set $K \subset \mathbb{R}^4$ on which:

- 1335 • the denominators m_1^2 and $j_2 - j_1^2$ in the definition of $\Phi(m_1, m_2, j_1, j_2)$ are bounded away
 - 1336 from zero by m_{\min}^2 and v_{\min} , and
 - 1337 • Φ is continuously differentiable.
- $$1338$$
- $$1339$$

1340 In particular, the gradient is bounded on K :

$$1341 \quad C_\Phi := \sup_{v \in K} \|\nabla \Phi(v)\| < \infty.$$

$$1342$$

$$1343$$

1344 Now, for any $v, \tilde{v} \in K$, consider the line segment $\gamma(t) := v + t(\tilde{v} - v)$, $t \in [0, 1]$, and define

1345 $\psi(t) := \Phi(\gamma(t))$. By the chain rule,

$$1346 \quad \psi'(t) = \nabla \Phi(\gamma(t)) \cdot (\tilde{v} - v),$$

$$1347$$

$$1348$$

1349 so that

$$|\psi'(t)| \leq \|\nabla \Phi(\gamma(t))\| \|\tilde{v} - v\| \leq C_\Phi \|\tilde{v} - v\| \quad \text{for all } t \in [0, 1].$$

1350 Integrating from 0 to 1, we obtain

$$\begin{aligned}
1351 & |\Phi(\tilde{v}) - \Phi(v)| = |\psi(1) - \psi(0)| \\
1352 & = \left| \int_0^1 \psi'(t) dt \right| \\
1353 & \leq \int_0^1 |\psi'(t)| dt \\
1354 & \leq \int_0^1 C_\Phi \|\tilde{v} - v\| dt \\
1355 & = C_\Phi \|\tilde{v} - v\|.
\end{aligned}$$

1362 Therefore

$$1363 \quad |\Phi(\tilde{v}) - \Phi(v)| \leq C_\Phi \|\tilde{v} - v\|. \quad (\text{A15})$$

1364 **Step 3: Combining the two bounds.** Finally, for h, \tilde{h} we have

$$1366 \quad |R(h)^2 - R(\tilde{h})^2| = |\Phi(v(h)) - \Phi(v(\tilde{h}))|.$$

1368 Applying equation A15 with $v = v(h)$ and $\tilde{v} = v(\tilde{h})$, and then equation A14, we obtain

$$1369 \quad |R(h)^2 - R(\tilde{h})^2| \leq C_\Phi \|v(h) - v(\tilde{h})\| \leq C_\Phi C_v \|h - \tilde{h}\|_\infty,$$

1371 where $K = C_\Phi C_v$. □

1373 **Lemma 2** (Approximation error of \tilde{p}_1 and \tilde{p}_n). *Let $I = [a, b]$ be the effective window with length $L = b - a$. Assume that p is twice continuously differentiable on I , and define*

$$1375 \quad C_{p''} := \sup_{u \in I} |p''(u)| < \infty.$$

1377 Then, for all $u \in I$,

$$1379 \quad |p(u) - \tilde{p}_1(u)| \leq \frac{1}{2} C_{p''} L^2, \quad |p(u) - \tilde{p}_n(u)| \leq \frac{1}{2} C_{p''} \frac{L^2}{n^2},$$

1381 and hence

$$1382 \quad \|\tilde{p}_1 - \tilde{p}_n\|_\infty \leq \frac{1}{2} C_{p''} L^2 \left(1 + \frac{1}{n^2}\right). \quad (\text{A16})$$

1385 **Error of the one-segment linear approximation \tilde{p}_1 .** Let $I = [a, b]$ and $L = b - a$. The one-segment linear model \tilde{p}_1 is defined as the first-order Taylor polynomial of p at $u = a$:

$$1387 \quad \tilde{p}_1(u) := p(a) + p'(a)(u - a), \quad u \in I.$$

1389 By Taylor's theorem with Lagrange remainder at $u = a$, for each $u \in I$ there exists a point ξ_u on the segment between a and u such that

$$1391 \quad p(u) = p(a) + p'(a)(u - a) + \frac{1}{2} p''(\xi_u)(u - a)^2.$$

1393 Hence

$$1394 \quad p(u) - \tilde{p}_1(u) = \frac{1}{2} p''(\xi_u)(u - a)^2.$$

1396 Taking absolute values and using the definition of $C_{p''}$,

$$1398 \quad |p(u) - \tilde{p}_1(u)| = \frac{1}{2} |p''(\xi_u)| |u - a|^2 \leq \frac{1}{2} C_{p''} |u - a|^2.$$

1400 Since $u \in [a, b]$ implies $|u - a| \leq L$, we obtain

$$1402 \quad |p(u) - \tilde{p}_1(u)| \leq \frac{1}{2} C_{p''} L^2 \quad \text{for all } u \in I.$$

1403

Error of the n -segment PWLA approximation \tilde{p}_n . Partition $I = [a, b]$ into n equal sub-intervals of length $\Delta := L/n$, and let $u_k := a + k\Delta$ for $k = 0, \dots, n$. On each sub-interval $[u_k, u_{k+1}]$, the n -segment PWLA model is

$$\tilde{p}_n(u) := p(u_k) + p'(u_k)(u - u_k), \quad u \in [u_k, u_{k+1}].$$

Applying Taylor's theorem at u_k for $u \in [u_k, u_{k+1}]$, there exists $\xi_{k,u} \in [u_k, u_{k+1}]$ such that

$$p(u) = p(u_k) + p'(u_k)(u - u_k) + \frac{1}{2}p''(\xi_{k,u})(u - u_k)^2.$$

Therefore

$$p(u) - \tilde{p}_n(u) = \frac{1}{2}p''(\xi_{k,u})(u - u_k)^2,$$

and thus

$$|p(u) - \tilde{p}_n(u)| = \frac{1}{2}|p''(\xi_{k,u})||u - u_k|^2 \leq \frac{1}{2}C_{p''}|u - u_k|^2.$$

Since $u \in [u_k, u_{k+1}]$ implies $|u - u_k| \leq \Delta = L/n$, we obtain

$$|p(u) - \tilde{p}_n(u)| \leq \frac{1}{2}C_{p''}\Delta^2 = \frac{1}{2}C_{p''}\frac{L^2}{n^2} \quad \text{for all } u \in I.$$

Sup-norm bound between \tilde{p}_1 and \tilde{p}_n . Taking the supremum over $u \in I$ in the pointwise bounds above gives

$$\|\tilde{p}_1 - p\|_\infty \leq \frac{1}{2}C_{p''}L^2, \quad \|p - \tilde{p}_n\|_\infty \leq \frac{1}{2}C_{p''}\frac{L^2}{n^2}.$$

By the triangle inequality in the sup norm,

$$\|\tilde{p}_1 - \tilde{p}_n\|_\infty \leq \|\tilde{p}_1 - p\|_\infty + \|p - \tilde{p}_n\|_\infty,$$

then,

$$\|\tilde{p}_1 - \tilde{p}_n\|_\infty \leq \|\tilde{p}_1 - p\|_\infty + \|p - \tilde{p}_n\|_\infty, \leq \frac{1}{2}C_{p''}L^2\left(1 + \frac{1}{n^2}\right).$$

□

Corollary A1 (Robustness of Theorem 2 under n -segment PWLA). *Suppose that, on the parameter range considered in this paper, the one-segment linear model \tilde{p}_1 used in Theorem 2 satisfies*

$$R(\tilde{p}_1)^2 \leq 1 - \delta_0$$

for some margin $\delta_0 > 0$ (for example, under the condition $L\delta > \sigma^2$ in Theorem 2). Let \tilde{p}_n be the n -segment piecewise-linear model defined above, and assume the conditions of Lemmas 1 and 2 hold. Then, for all $n \in \mathbb{N}$,

$$|R(\tilde{p}_n)^2 - R(\tilde{p}_1)^2| \leq K \|\tilde{p}_n - \tilde{p}_1\|_\infty \leq \frac{1}{2}KC_{p''}L^2\left(1 + \frac{1}{n^2}\right),$$

where $K > 0$ is the Lipschitz constant from Lemma 1 (measuring the sensitivity of $R(h)^2$ to perturbations of the weight h) and $C_{p''} := \sup_{u \in I} |p''(u)|$ is the curvature bound from Lemma 2 (quantifying how strongly the Gaussian weight can bend on the effective window). In particular, if the effective window satisfies

$$KC_{p''}L^2 \leq \delta_0,$$

then

$$R(\tilde{p}_n)^2 \leq R(\tilde{p}_1)^2 + \frac{\delta_0}{2} \leq 1 - \frac{\delta_0}{2} < 1.$$

Hence, the inequality $\text{CV}_{\text{asy}} < \text{CV}_{\text{sym}}$ derived under the one-segment linear model in Theorem 2 can be valid for any sufficiently accurate n -segment piecewise-linear approximation \tilde{p}_n of the same Gaussian weight function on the effective window.

1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511

Table A3: Results of accuracy and spike count under different deletion rates on CIFAR10 with VGG16.

| Methods | Ratio(%) | Acc.(%) | # of Spikes ($\times 10^3$) |
|----------------------------------|----------|-----------------------|----------------------------------|
| <i>BOX</i> (Baseline) | 0 | 94.30 (-0.00) | 94 (-0.0%) |
| | 10 | 91.09 (-3.21) | 88 (-6.0%) |
| | 20 | 69.08 (-25.22) | 83 (-14.0%) |
| <i>A²SG</i> (Ours) | 0 | 95.29 (-0.00) | 89 (-0.0%) |
| | 10 | 92.59 (-2.81) | 84 (-5.0%) |
| | 20 | 69.08 (-24.89) | 81 (-10.0%) |