# CARE-BN: PRECISE MOVING STATISTICS FOR STABILIZING SPIKING NEURAL NETWORKS IN REINFORCEMENT LEARNING

Zijie Xu, Xinyu Shi, Yiting Dong, Zihan Huang, Zhaofei Yu<sup>†</sup> Peking University Beijing, 100871, China

### **ABSTRACT**

Spiking Neural Networks (SNNs) offer low-latency and energy-efficient decisionmaking on neuromorphic hardware by mimicking the event-driven dynamics of biological neurons. However, due to the discrete and non-differentiable nature of spikes, directly trained SNNs rely heavily on Batch Normalization (BN) to stabilize gradient updates. In online Reinforcement Learning (RL), imprecise BN statistics hinder exploitation, resulting in slower convergence and suboptimal policies. This challenge limits the adoption of SNNs for energy-efficient control on resource-constrained devices. To overcome this, we propose Confidenceadaptive and Re-calibration Batch Normalization (CaRe-BN), which introduces (i) a confidence-guided adaptive update strategy for BN statistics and (ii) a recalibration mechanism to align distributions. By providing more accurate normalization, CaRe-BN stabilizes SNN optimization without disrupting the RL training process. Importantly, CaRe-BN does not alter inference, thus preserving the energy efficiency of SNNs in deployment. Extensive experiments on continuous control benchmarks demonstrate that CaRe-BN improves SNN performance by up to 22.6% across different spiking neuron models and RL algorithms. Remarkably, SNNs equipped with CaRe-BN even surpass their ANN counterparts by 5.9%. These results highlight a new direction for BN techniques tailored to RL, paving the way for neuromorphic agents that are both efficient and high-performing.

### 1 Introduction

Spiking Neural Networks (SNNs) have emerged as a promising class of neural models that more closely mimic the event-driven computation of biological brains (Maass, 1997; Gerstner et al., 2014). This event-driven property makes SNNs particularly well suited for deployment on neuromorphic hardware platforms (Davies et al., 2018; DeBole et al., 2019), enabling low-latency and energy-efficient inference.

In parallel, Reinforcement Learning (RL) has achieved remarkable success across a wide range of domains (Mnih et al., 2015b; Lillicrap et al., 2015; Haarnoja et al., 2018). Among these, continuous control tasks have received significant attention due to their alignment with real-world scenarios and their strong connection to embodied AI and robotic applications (Kober et al., 2013; Gu et al., 2017; Brunke et al., 2022). Combining the strengths of SNNs with RL (SNN-RL) offers the potential to train agents that not only learn complex behaviors but also execute them with extremely low energy consumption (Yamazaki et al., 2022). This makes SNN-RL particularly appealing for robotics and autonomous systems deployed on resource-constrained edge devices.

However, training SNNs is challenging. Due to their discrete dynamics and non-differentiable activation functions, SNNs are heavily dependent on Batch Normalization (BN) (Ioffe & Szegedy, 2015) to mitigate gradient vanishing or explosion. BN stabilizes optimization by normalizing activations with moving estimates of the mean and variance, thereby alleviating temporal covariate shifts in SNNs and enabling state-of-the-art performance (Duan et al., 2022; Jiang et al., 2024).

While effective in supervised learning, BN suffers a severe breakdown in online RL because moving statistics cannot be estimated precisely under nonstationary dynamics. As shown in Figure 1,

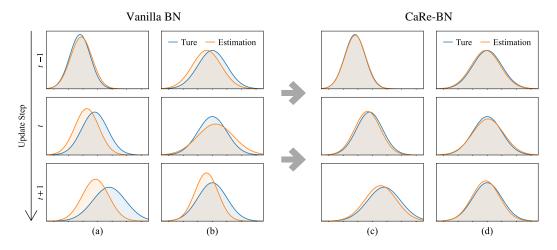


Figure 1: Real and estimated input activation distributions in BN layers. Distributions change rapidly in (a) and (c), while remaining stable in (b) and (d).

traditional BN struggles to track the true statistics: When distributions shift rapidly (Figure 1(a)), estimates lag behind; when distributions are relatively static (Figure 1(b)), estimates contain noise. These inaccuracies lead agents to select suboptimal actions and generate poor trajectories, which are then reused for training—further compounding the problem and hindering policy improvement.

In this work, we address this issue by proposing **Confidence-adaptive and Re-calibration Batch Normalization** (CaRe-BN), a BN strategy tailored for SNN-based RL. CaRe-BN introduces two complementary components: (*i*) **Confidence-adaptive update** (Ca-BN), a confidence-weighted moving estimator of BN statistics that ensures unbiasedness and optimal variance reduction; and (*ii*) **Re-calibration** (Re-BN), a periodic correction scheme that leverages replay buffer resampling to refine inference statistics. Together, these mechanisms enable precise, low-variance estimation of BN statistics under the nonstationary dynamics of SNN-RL (Figure 1). With more accurate moving statistics, CaRe-BN stabilizes SNN optimization **without disrupting the online RL process**.

We evaluate CaRe-BN on various continuous control tasks from MuJoCo (Todorov et al., 2012; Todorov, 2014b). The results show that CaRe-BN not only resolves the issue of imprecise BN statistics but also accelerates training and achieves state-of-the-art performance. Remarkably, SNN-based agents equipped with CaRe-BN even **outperform their ANN counterparts by 5.9**%, without requiring complex neuron dynamics or specialized RL frameworks.

# 2 RELATED WORKS

# 2.1 BATCH NORMALIZATION IN SPIKING NEURAL NETWORKS

Batch Normalization (BN) was originally proposed for ANNs to mitigate internal covariate shift during training (Ioffe & Szegedy, 2015), thereby accelerating convergence and improving performance (Santurkar et al., 2018). To address unstable training in SNNs, several extensions of BN have been developed (Zheng et al., 2021; Duan et al., 2022; Kim & Panda, 2021; Jiang et al., 2024). While these methods are effective in supervised tasks, they are designed under the assumption of static training—inference distributions. This assumption is violated in online RL, where distributions shift continually as the agent interacts with the environment, making these BN variants ill-suited for SNN-RL.

### 2.2 Spiking Neural Networks in Reinforcement Learning

Early work in SNN-RL primarily relied on synaptic plasticity rules, particularly reward-modulated Spike-Timing-Dependent Plasticity (R-STDP) and its variants (Florian, 2007; Frémaux & Gerstner, 2016; Gerstner et al., 2018; Frémaux et al., 2013; Yang et al., 2024). Another research direction

focused on ANN-to-SNN conversion: Patel et al. (2019); Tan et al. (2021); Kumar et al. (2025) converted Deep Q-Networks (DQNs) (Mnih, 2013; Mnih et al., 2015a) into SNNs. To enable direct gradient-based training, Liu et al. (2022); Chen et al. (2022); Qin et al. (2022); Sun et al. (2022) applied Spatio-Temporal Backpropagation (STBP) (Wu et al., 2018) to train DQNs, while Bellec et al. (2020) introduced e-prop with eligibility traces to train policy networks using policy gradient methods (Sutton et al., 1999).

For continuous control tasks, hybrid frameworks have been extensively explored (Tang et al., 2020; 2021; Zhang et al., 2022; Chen et al., 2024a; Zhang et al., 2024; Ding et al., 2022; Chen et al., 2024b; Xu et al., 2025). These approaches typically employ a Spiking Actor Network (SAN) cotrained with a deep ANN critic in the Actor–Critic framework (Konda & Tsitsiklis, 1999). However, none of these methods address the challenge of normalization in SNN-based RL. The absence of proper normalization often leads to unstable updates, slower convergence, or even divergence during training.

### 3 PRELIMINARIES

# 3.1 Spiking Neural Networks

Spiking Neural Networks (SNNs) communicate through discrete spikes rather than continuous activations. The most widely used neuron model is the Leaky Integrate-and-Fire (LIF) neuron, whose membrane potential dynamics are described as:

$$H_t = \lambda V_{t-1} + C_t, \qquad S_t = \Theta(H_t - V_{th}), \qquad V_t = (1 - S_t) \cdot H_t + S_t \cdot V_{\text{reset}}, \tag{1}$$

where  $C_t$ ,  $H_t$ ,  $S_t$ , and  $V_t$  denote the input current, the accumulated membrane potential, the binary output spike, and the post-firing membrane potential at time step t, respectively. The parameters  $V_{th}$ ,  $V_{\text{reset}}$ , and  $\lambda$  represent the firing threshold, reset voltage, and leakage factor, respectively.  $\Theta(\cdot)$  is the Heaviside step function.

### 3.2 Reinforcement Learning

Reinforcement Learning (RL) is a framework in which an agent learns to maximize cumulative rewards by interacting with an environment. The agent maps states (or observations) to actions, with the learning loop consisting of two steps: (i) the agent selects an action, receives a reward, and transitions to the next state; and (ii) the agent updates its policy by sampling mini-batches of past experiences.

Because the policy continuously evolves during training, the data distribution is inherently non-stationary. This poses challenges for batch normalization methods, which rely on the assumption of a stationary distribution.

### 3.3 BATCH NORMALIZATION

Batch Normalization (BN) (Ioffe & Szegedy, 2015) is a widely used technique to stabilize and accelerate the training of deep neural networks. Given an activation  $x_i \in \mathbb{R}^d$  at iteration i, BN normalizes it using the mean and variance computed over a mini-batch  $\mathcal{B} = \{x_i^1, \dots, x_i^N\}$ :

$$\mu_{\mathcal{B}} = \frac{1}{N} \sum_{j=1}^{N} x_i^j, \quad \sigma_{\mathcal{B}}^2 = \frac{1}{N} \sum_{j=1}^{N} (x_i^j - \mu_{\mathcal{B}})^2,$$
 (2)

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}, \quad y_i = \gamma \hat{x}_i + \beta, \tag{3}$$

where  $\epsilon$  is a small constant for numerical stability, and  $\gamma, \beta$  are learnable affine parameters. During inference, moving statistics  $(\hat{\mu}_i, \hat{\sigma}_i^2)$  are used in place of batch statistics  $(\mu_i, \sigma_i^2)$ .

In supervised learning, this discrepancy between training (mini-batch statistics) and inference (moving statistics) is usually tolerable, as imprecise moving estimates do not directly affect gradient updates. However, in online RL, inaccurate moving statistics degrade policy exploitation, leading to unstable training dynamics and even divergence.

# 4 METHODOLOGY

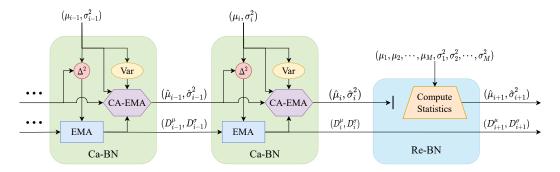


Figure 2: The statistics estimation scheme of CaRe-BN. In this framework, Ca-BN is applied at every update step, while Re-BN is performed periodically.  $\Delta^2$  denotes the squared error, **Var** represents the variance computed according to Eq. 9, **EMA** refers to the exponential moving average in Eq. 12, and **CA-EMA** denotes the confidence-adaptive update defined in Eqs. 5 and 6.

As illustrated in Figure 2, we propose Confidence-adaptive and Recalibration Batch Normalization (CaRe-BN) to address the challenge of approximating moving statistics in online RL. Section 4.1 analyzes the limitations of traditional BN in online RL, where statistics are often estimated imprecisely. Section 4.2 introduces the confidence-adaptive update mechanism (Ca-BN), which dynamically adjusts statistics estimation based on the reliability of the current approximation. Section 4.3 presents the recalibration mechanism (Re-BN), which periodically corrects accumulated estimation errors. Finally, Section 4.4 integrates these components into the full CaRe-BN framework and demonstrates its use in online RL algorithms.

### 4.1 ISSUES IN APPROXIMATING MOVING STATISTICS

Online RL introduces stronger distribution shifts. Unlike supervised learning, where the data distribution is typically assumed to be static, online RL involves continuous interaction between the agent and the environment. This results in a non-stationary data distribution, which in turn causes activation statistics to drift over time.

**Inaccurate statistics degrade RL performance.** Supervised learning only requires the final moving statistics to be accurate, as inference is performed after training. In contrast, online RL requires reliable statistics throughout training. When statistics are imprecise, the agent selects suboptimal actions during exploration and exploitation, generating poor trajectories that further degrades policy updates.

The key of the problem lies in accurately estimating inference-time statistics under shifting distributions. Hence, it is essential to design estimators that adapt to distributional changes while minimizing approximation error during training.

It is worth noting that most conventional ANN-based RL algorithms do not employ BN (Lillicrap, 2015; Sutton & Barto, 2018), as shallow ANNs can often learn stable representations without normalization. In contrast, BN is indispensable for stabilizing SNNs training. Therefore, addressing this issue is particularly critical for SNN-based RL.

# 4.2 CONFIDENCE-ADAPTIVE UPDATE OF BN STATISTICS (CA-BN)

Conventional BN approximates population statistics using an exponential moving average (EMA) of the batch mean and variance:

$$\hat{\mu}_i \leftarrow (1 - \alpha)\hat{\mu}_{i-1} + \alpha\mu_i, \quad \hat{\sigma}_i^2 \leftarrow (1 - \alpha)\hat{\sigma}_{i-1}^2 + \alpha\sigma_i^2, \tag{4}$$

where  $\alpha$  is the momentum parameter. This update rule faces a fundamental **noise–delay trade-off**. As shown in Figure 1, low momentum yields stable but slow adaptation to distribution shifts, while

high momentum adapts quickly but amplifies the noise from small-batch estimates. This trade-off is particularly harmful in online RL, where accurate normalization is critical for stable policy learning.

Inspired by the Kalman estimator (Kalman, 1960), we derive a confidence-guided mechanism that adaptively reweights estimators to minimize the mean-squared error (MSE) of BN statistics.

**Theorem 1** Let  $(\mu_i, \sigma_i^2)$  and  $(\hat{\mu}_{i|i-1}, \hat{\sigma}_{i|i-1}^2)$  be two unbiased estimators of the population parameters  $(\mu_i^*, \sigma_i^{*2})$ . The optimal linear estimator that minimizes MSE is

$$\hat{\mu}_i = (1 - K_i^{\mu})\hat{\mu}_{i|i-1} + K_i^{\mu}\mu_i, \qquad K_i^{\mu} = \frac{\mathbb{D}(\mu_i^* - \hat{\mu}_{i|i-1})}{\mathbb{D}(\mu_i^* - \hat{\mu}_{i|i-1}) + \mathbb{D}(\mu_i^* - \mu_i)},$$
(5)

$$\hat{\sigma}_{i}^{2} = (1 - K_{i}^{\sigma})\hat{\sigma}_{i|i-1}^{2} + K_{i}^{\sigma}\sigma_{i}^{2}, \qquad K_{i}^{\sigma} = \frac{\mathbb{D}(\sigma_{i}^{*2} - \hat{\sigma}_{i|i-1}^{2})}{\mathbb{D}(\sigma_{i}^{*2} - \hat{\sigma}_{i|i-1}^{2}) + \mathbb{D}(\sigma_{i}^{*2} - \sigma_{i}^{2})}, \tag{6}$$

where  $K_i^{\mu}$  and  $K_i^{\sigma}$  are confidence-guided adaptive weights, and  $\mathbb{D}(\cdot)$  denotes estimation variance (the reciprocal of confidence).

**Proof 1** Since both  $\hat{\mu}_{i|i-1}$  and  $\mu_i$  are unbiased for  $\mu_i^*$ , any linear combination  $\tilde{\mu}_i = (1-K)\hat{\mu}_{i|i-1} + K\mu_i$  is also unbiased. The variance is

$$\mathbb{D}(\tilde{\mu}_i - \mu_i^*) = (1 - K)^2 \cdot \mathbb{D}(\hat{\mu}_{i|i-1} - \mu_i^*) + K^2 \cdot \mathbb{D}(\mu_i - \mu_i^*). \tag{7}$$

Minimizing over K yields the optimal  $K = K_i^{\mu}$ . The variance update (Eq. 6) follows analogously.

**Assumption 1** The activations in iteration i are modeled as  $x_i \sim \mathcal{N}(\mu_i^*, \sigma_i^{*2})$ , following the standard Gaussianity assumption in BN.

Confidence of mini-batch statistics. For a batch of size N, the sample mean  $\mu_i$  and variance  $\sigma_i^2$  satisfy

$$\mu_i \sim \mathcal{N}\left(\mu_i^*, \frac{\sigma_i^{*2}}{N}\right), \qquad \frac{(N-1)\sigma_i^2}{\sigma_i^{*2}} \sim \chi_{N-1}^2.$$
 (8)

Since  $\mu_i$  and  $\sigma_i^2$  are unbiased estimators of  $\mu_i^*$  and  $\sigma_i^{*2}$ , we approximate

$$\mathbb{D}(\mu_i^* - \mu_i) = \frac{\sigma_i^{*2}}{N} \approx \frac{\sigma_i^2}{N}, \qquad \mathbb{D}(\sigma_i^{*2} - \sigma_i^2) = \frac{2\sigma_i^{*4}}{N - 1} \approx \frac{2\sigma_i^4}{N - 1}.$$
 (9)

Confidence of previous estimates. Since the true distribution drift is unknown, we approximate

$$\hat{\mu}_{i|i-1} = \hat{\mu}_{i-1}, \qquad \hat{\sigma}_{i|i-1}^2 = \hat{\sigma}_{i-1}^2,$$
 (10)

$$\mathbb{D}(\mu_i^* - \hat{\mu}_{i|i-1}) \approx D_i^{\mu}, \qquad \mathbb{D}(\sigma_i^{*2} - \hat{\sigma}_{i|i-1}^2) \approx D_i^{\sigma}, \tag{11}$$

where  $D_i^{\mu}$ ,  $D_i^{\sigma}$  are recursively smoothed statistics with momentum  $\alpha$ :

$$D_i^{\mu} \leftarrow (1 - \alpha) D_{i-1}^{\mu} + \alpha (\mu_i - \hat{\mu}_{i-1})^2, \qquad D_i^{\sigma} \leftarrow (1 - \alpha) D_{i-1}^{\sigma} + \alpha (\sigma_i^2 - \hat{\sigma}_{i-1}^2)^2.$$
 (12)

Combining Eqs. 5–12, we obtain the confidence-adaptive update scheme. When distributional shifts are rapid,  $D_i^\mu$  and  $D_i^\sigma$  grow large, increasing  $K_i^\mu$  and  $K_i^\sigma$  and accelerating adaptation. Conversely, when statistics are stable, these terms shrink, lowering  $K_i^\mu$  and  $K_i^\sigma$  and reducing noise from small mini-batches.

### 4.3 RE-CALIBRATION MECHANISM OF BN STATISTICS (RE-BN)

While the confidence-adaptive update provides online estimates of BN statistics during training, these estimates may still drift from the true population values due to stochastic mini-batch noise. The most accurate approach would be to recompute exact statistics by forward-propagating the entire dataset after each update (Wu & Johnson, 2021). However, this is computationally infeasible in RL, as it would require processing millions of samples at every step.

A more practical alternative is to periodically re-calibrate BN statistics using larger aggregated batches. Specifically, at fixed intervals  $T_{\text{cal}}$ , we draw M calibration batches  $\{\mathcal{B}_1,\ldots,\mathcal{B}_M\}$  from the replay buffer. For each batch  $\mathcal{B}_j$ , we compute its mean  $\mu_j$  and variance  $\sigma_j^2$ . The recalibrated BN statistics are then given by:

$$\hat{\mu}_i = \frac{1}{M} \sum_{j=1}^M \mu_j, \qquad \hat{\sigma}_i^2 = \frac{1}{M} \sum_{j=1}^M (\sigma_j^2 + \mu_j^2) - \hat{\mu}_i^2.$$
 (13)

This recalibration requires additional forward passes, but the extra overhead is upper bounded by  $\frac{M}{T_{\rm cal}}$  times the total training cost. Since we set  $T_{\rm cal} \gg M$ , the computational overhead remains negligible, while significantly improving the accuracy of BN statistics.

### 4.4 INTEGRATING WITH RL

The proposed Confidence-adaptive and Re-calibration Batch Normalization (CaRe-BN) integrates two complementary mechanisms: the confidence-adaptive update in Section 4.2, which provides an online estimation of batch normalization (BN) statistics, and the re-calibration procedure in Section 4.3, which corrects accumulated bias. The overall integration within an online RL framework is outlined in Algorithm 1.

# Algorithm 1 General RL Algorithm with CaRe-BN

- 1: Initialize the agent networks and the replay buffer.
- 2: for each iteration do
- 3: Select an action and store the transition (**inference BN statistics**).
- 4: Update the agent by sampling a minibatch of N transitions (mini-batch BN statistics).
- 5: Update the moving BN statistics as:

$$D_{i}^{\mu} \leftarrow (1 - \alpha)D_{i-1}^{\mu} + \alpha(\mu_{i} - \hat{\mu}_{i-1})^{2}, \qquad D_{i}^{\sigma} \leftarrow (1 - \alpha)D_{i-1}^{\sigma} + \alpha(\sigma_{i}^{2} - \hat{\sigma}_{i-1}^{2})^{2},$$
$$\hat{\mu}_{i} = \frac{D_{i}^{\mu} \cdot \mu_{i} + \frac{\sigma_{i}^{2}}{N} \cdot \hat{\mu}_{i-1}}{D_{i}^{\mu} + \frac{\sigma_{i}^{2}}{N}}, \qquad \hat{\sigma}_{i}^{2} = \frac{D_{i}^{\sigma} \cdot \sigma_{i}^{2} + \frac{2\sigma_{i}^{4}}{N-1} \cdot \hat{\sigma}_{i-1}^{2}}{D_{i}^{\sigma} + \frac{2\sigma_{i}^{4}}{N-1}}.$$

- 6: **if** Re-calibration **then**
- 7: Sample M minibatches of N transitions each and update BN statistics using Eq. (13).
- 8: end if
- 9: end for

It is important to note that the inference procedure of CaRe-BN remains identical to that of conventional BN. Consequently, the CaRe-BN layer is seamlessly fused into synaptic weights, introducing no additional inference overhead during deployment.

### 5 EXPERIMENTS

# 5.1 EXPERIMENT SETUP

We evaluate our CaRe-BN on five standard continuous control benchmarks from the MuJoCo suite (Todorov et al., 2012; Todorov, 2014b) within the OpenAI Gymnasium (Brockman, 2016; Towers et al., 2024), including InvertedDoublePendulum (IDP) (Todorov, 2014a), Ant (Schulman et al., 2015), HalfCheetah (Wawrzyński, 2009), Hopper (Erez et al., 2012), and Walker2d. All environments use default settings, and performance is evaluated by averaging the rewards over 10 trials.

To assess the effectiveness and adaptability of CaRe-BN, we conduct experiments with different spiking neuron models and RL algorithms. For spiking neurons, we test the Leaky Integrate-and-Fire (LIF) (Gerstner & Kistler, 2002) and Current-based LIF (CLIF) (Tang et al., 2021), with the dynamics detailed in Appendix. For RL algorithms, we employ Deep Deterministic Policy Gradient (DDPG) (Lillicrap, 2015) and Twin Delayed DDPG (TD3) (Fujimoto et al., 2018). All SNNs are simulated with a time step of 5, and additional training details are provided in Appendix.

For a fair comparison, we maintain consistent hyperparameters across all models and algorithms, with full settings provided in Appendix. Each algorithm is trained using 5 random seeds, and we report the averaged results.

### 5.2 More Precise BN Statistics Lead to Better Exploration

A key aspect of online RL is that the quality of exploration significantly influences subsequent policy updates. Figure 3 compares the exploration performance of standard BN with our proposed CaRe-BN. Throughout the training process, CaRe-BN consistently achieves higher exploration returns.

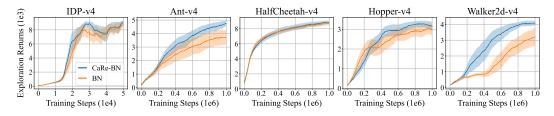


Figure 3: Exploration performance of BN and CaRe-BN with CLIF neurons and the TD3 algorithm across five MuJoCo tasks. Shaded areas represent half a standard deviation across five random seeds. Curves are uniformly smoothed for visual clarity.

Since CaRe-BN does not directly modify the gradient update process, the observed improvement in exploration performance is solely due to its more precise estimation of BN statistics. This leads to better exploration policies, which in turn generate higher-quality trajectories for updating the agent. As a result, CaRe-BN forms a positive feedback loop: improved statistics  $\rightarrow$  better exploration  $\rightarrow$  higher-quality experiences  $\rightarrow$  better policy.

### 5.3 ADAPTABILITY OF CARE-BN

To evaluate the adaptability of CaRe-BN, we test it across different RL algorithms (i.e., DDPG (Lillicrap, 2015) and TD3 (Fujimoto et al., 2018)) and spiking neuron models (i.e., LIF and CLIF (Tang et al., 2021)).

**Better final return.** Figure 4 shows the learning curves for SNN models with and without CaRe-BN. In most cases, CaRe-BN consistently outperforms standard SNNs, converging faster and achieving higher final returns. These improvements are robust across different spiking neurons and RL algorithms, confirming that CaRe-BN enhances performance in diverse settings.

**Lower variance.** Figure 5 (a) and (b) display the relative variance of the final policy. Compared to standard SNNs, CaRe-BN significantly reduces the variance of SNN-RL training, and even achieves lower variance than ANN baselines (i.e., 17.71% for DDPG and 21.24% for TD3). This indicates that CaRe-BN not only enhances performance but also improves the stability and reproducibility of SNN-RL.

### 5.4 EXCEEDING SOTA

To further validate the effectiveness of CaRe-BN, we compare it with existing state-of-the-art (SOTA) SNN-RL methods and various batch normalization strategies for SNNs. The evaluation is conducted using the TD3 algorithm (Fujimoto et al., 2018) (a strong SOTA baseline for continuous control) and the CLIF neuron model (Tang et al., 2021) (the most commonly used neuron type in recent SNN-RL studies). The ANN-SNN conversion baseline follows the SOTA method proposed in Bu et al. (2025). For direct-trained SNNs, we include pop-SAN (Tang et al., 2021), MDC-SAN (Zhang et al., 2022), and ILC-SAN (Chen et al., 2024a). Additionally, we test several BN algorithms for SNNs, including tdBN (Zheng et al., 2021), BNTT (Kim & Panda, 2021), TEBN (Duan et al., 2022), and TABN (Jiang et al., 2024). The performance is summarized in Table 1, where the average

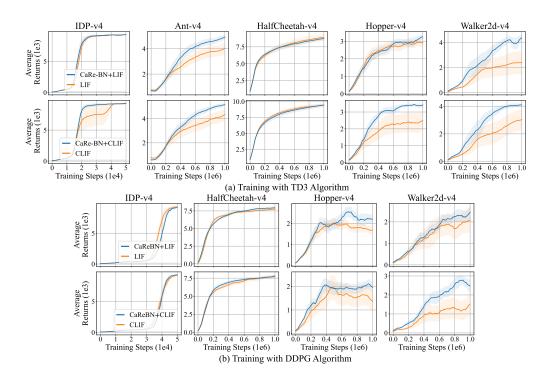


Figure 4: Learning curves of SNN-based agents trained with TD3 (top) and DDPG (bottom). Since the DDPG algorithm (in both ANN and SNN) diverges in the Ant-v4 environment, these curves are not shown. Shaded areas represent half a standard deviation across five random seeds. Curves are uniformly smoothed for visual clarity.

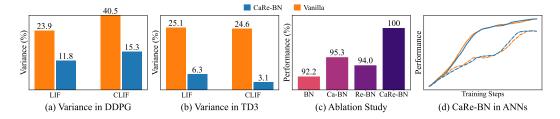


Figure 5: (a), (b) Relative variance percentage of final policy returns, computed by averaging the standard deviation ratio across five random seeds, for all environments. (c) Normalized maximum performance across all environments for the ablation study. (d) Normalized learning curves across all environments for ANNs implementing CaRe-BN. The dashed lines represent DDPG and the solid lines represent TD3. Performance and training steps are normalized linearly. Curves are uniformly smoothed for visual clarity.

performance gain (APG) is defined as:

$$APG = \left(\frac{1}{|\text{envs}|} \sum_{\text{env} \in \text{envs}} \frac{\text{performance}(\text{env})}{\text{baseline}(\text{env})} - 1\right) \cdot 100\%,\tag{14}$$

where |envs| denotes the total number of environments, and performance(env) and baseline(env) represent the performance of the evaluated algorithm and the ANN baseline in each environment, respectively.

**Compared with other SNN-RL methods:** CaRe-BN significantly outperforms previous SNN-RL approaches, demonstrating that normalization plays a more crucial role than architectural modifications in improving SNN-RL performance.

**Compared with other BN methods:** Compared to existing SNN-specific BN variants, CaRe-BN performs superior, establishing a new state-of-the-art normalization strategy for SNN-RL.

Compared with ANNs: Notably, CaRe-BN outperforms the ANN baseline by 5.9% on average. This highlights that with proper normalization, SNNs can not only match but exceed the performance of traditional ANN-based RL agents, while retaining their energy-efficient advantages.

Table 1: Max average returns over 5 random seeds with CLIF spiking neurons, and the average performance gain (APG) against ANN baseline, where  $\pm$  denotes one standard deviation.

Method	IDP-v4	Ant-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4	APG
ANN (TD3)	$7503 \pm 3713$	$4770 \pm 1014$	$10857 \pm 475$	$3410 \pm 164$	$4340 \pm 383$	0.00%
ANN-SNN	$3859 \pm 4440$	$3550 \pm 963$	$8703 \pm 658$	$3098 \pm 281$	$4235 \pm 354$	-21.11%
pop-SAN	$9351 \pm 1$	$4590 \pm 1006$	$9594 \pm 689$	$2772 \pm 1263$	$3307 \pm 1514$	-6.66%
MDC-SAN	$9350 \pm 1$	$4800 \pm 994$	$9147 \pm 231$	$3446 \pm 131$	$3964 \pm 1353$	0.37%
ILC-SAN	$9352\pm1$	$5584 \pm 272$	$9222 \pm 615$	$3403\pm148$	$4200 \pm 717$	4.64%
tdBN	$9346 \pm 2$	$4403 \pm 1134$	$9402 \pm 527$	$3592 \pm 46$	$3464 \pm 970$	-2.28%
BNTT	$9347 \pm 1$	$4379 \pm 941$	$9466 \pm 659$	$3524 \pm 161$	$3689 \pm 1247$	-1.62%
TEBN	$9349 \pm 1$	$4408 \pm 1156$	$9452 \pm 539$	$3472 \pm 135$	$4235 \pm 381$	0.69%
TABN	$9348 \pm 2$	$4382 \pm 753$	$9784 \pm 169$	$3585 \pm 83$	$4537 \pm 398$	3.25%
CaRe-BN	$9348 \pm 2$	$5373 \pm 159$	$9563 \pm 442$	$3586 \pm 49$	$4296 \pm 268$	<b>5.90</b> %

### 5.5 ABLATION STUDIES

We conduct ablation studies by separately evaluating the effects of the Confidence-adaptive update (Ca-BN) and the Re-calibration mechanism (Re-BN), as shown in Figure 5 (c). The results demonstrate that both the adaptive estimation and recalibration mechanisms are beneficial on their own. However, their combination provides the most significant improvement. Specifically, Ca-BN addresses the mismatch between training and inference statistics, while Re-BN corrects accumulated errors, further stabilizing training. By integrating both components, CaRe-BN achieves more precise and consistent normalization, leading to superior overall performance.

### 5.6 SNN-FRIENDLY DESIGN

Dispite the stunning improvement in SNNs, we also evaluate CaRe-BN on standard ANNs trained with TD3 and DDPG, as shown in Figure 5 (d). The results indicate that ANNs with CaRe-BN perform similarly to their baseline counterparts without CaRe-BN. This outcome is expected for the following reasons: (i) Shallow ANNs can already train stably and effectively without normalization<sup>1</sup>, so adding CaRe-BN does not provide significant improvements. (ii) While CaRe-BN provides more precise estimates of BN statistics, this does not negatively impact the RL training process. These results further underscore that the improvements observed are not due to a stronger RL mechanism, but rather to the SNN-specific normalization strategies.

### 6 CONCLUSION

In this work, we introduced CaRe-BN, the first batch normalization method specifically designed for SNNs in RL. By addressing the instability of conventional BN in online RL, CaRe-BN enables SNNs to outperform their ANN counterparts in continuous control tasks. Importantly, CaRe-BN is lightweight and easy to integrate, making it a seamless drop-in replacement for existing SNN-RL pipelines without introducing additional computational overhead.

Beyond its technical contributions, CaRe-BN brings SNN-RL one step closer to practical deployment. By stabilizing training and improving exploration, it unlocks the potential of SNNs to act as both energy-efficient and high-performance agents in real-world continuous control applications. We believe this work underscores the importance of normalization strategies tailored to the unique dynamics of SNNs and opens new avenues for bridging the gap between neuromorphic learning and reinforcement learning at scale.

<sup>&</sup>lt;sup>1</sup>In RL, networks typically consist of two hidden layers with 256 neurons

# REFERENCES

- Guillaume Bellec, Franz Scherr, Anand Subramoney, Elias Hajek, Darjan Salaj, Robert Legenstein, and Wolfgang Maass. A solution to the learning dilemma for recurrent networks of spiking neurons. *Nature communications*, 11(1):3625, 2020.
- G Brockman. Openai gym. arXiv preprint arXiv:1606.01540, 2016.
- Lukas Brunke, Melissa Greeff, Adam W Hall, Zhaocong Yuan, Siqi Zhou, Jacopo Panerati, and Angela P Schoellig. Safe learning in robotics: From learning-based control to safe reinforcement learning. *Annual Review of Control, Robotics, and Autonomous Systems*, 5(1):411–444, 2022.
- Tong Bu, Maohua Li, and Zhaofei Yu. Inference-scale complexity in ann-snn conversion for high-performance and low-power applications. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 24387–24397, 2025.
- Ding Chen, Peixi Peng, Tiejun Huang, and Yonghong Tian. Deep reinforcement learning with spiking q-learning. *arXiv preprint arXiv:2201.09754*, 2022.
- Ding Chen, Peixi Peng, Tiejun Huang, and Yonghong Tian. Fully spiking actor network with intralayer connections for reinforcement learning. *IEEE Transactions on Neural Networks and Learning Systems*, 36(2):2881–2893, 2024a.
- Ding Chen, Peixi Peng, Tiejun Huang, and Yonghong Tian. Noisy spiking actor network for exploration. *arXiv preprint arXiv:2403.04162*, 2024b.
- Mike Davies, Narayan Srinivasa, Tsung-Han Lin, Gautham Chinya, Yongqiang Cao, Sri Harsha Choday, Georgios Dimou, Prasad Joshi, Nabil Imam, Shweta Jain, et al. Loihi: A neuromorphic manycore processor with on-chip learning. *IEEE Micro*, 2018.
- Michael V DeBole, Brian Taba, Arnon Amir, Filipp Akopyan, Alexander Andreopoulos, William P Risk, Jeff Kusnitz, Carlos Ortega Otero, Tapan K Nayak, Rathinakumar Appuswamy, et al. TrueNorth: Accelerating from zero to 64 million neurons in 10 years. *Computer*, 2019.
- Jianchuan Ding, Bo Dong, Felix Heide, Yufei Ding, Yunduo Zhou, Baocai Yin, and Xin Yang. Biologically inspired dynamic thresholds for spiking neural networks. Advances in neural information processing systems, 35:6090–6103, 2022.
- Chaoteng Duan, Jianhao Ding, Shiyan Chen, Zhaofei Yu, and Tiejun Huang. Temporal effective batch normalization in spiking neural networks. *Advances in Neural Information Processing Systems*, 35:34377–34390, 2022.
- Tom Erez, Yuval Tassa, and Emanuel Todorov. Infinite-horizon model predictive control for periodic tasks with contacts. *Robotics: Science and Systems VII*, 2012.
- Răzvan V Florian. Reinforcement learning through modulation of spike-timing-dependent synaptic plasticity. *Neural computation*, 19(6):1468–1502, 2007.
- Nicolas Frémaux and Wulfram Gerstner. Neuromodulated spike-timing-dependent plasticity, and theory of three-factor learning rules. *Frontiers in neural circuits*, 9:85, 2016.
- Nicolas Frémaux, Henning Sprekeler, and Wulfram Gerstner. Reinforcement learning using a continuous time actor-critic framework with spiking neurons. *PLoS computational biology*, 9(4): e1003024, 2013.
- Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, pp. 1587–1596. PMLR, 2018.
- Wulfram Gerstner and Werner M Kistler. *Spiking neuron models: Single neurons, populations, plasticity.* Cambridge university press, 2002.
- Wulfram Gerstner, Werner M Kistler, Richard Naud, and Liam Paninski. *Neuronal dynamics: From single neurons to networks and models of cognition*. Cambridge University Press, 2014.

- Wulfram Gerstner, Marco Lehmann, Vasiliki Liakoni, Dane Corneil, and Johanni Brea. Eligibility traces and plasticity on behavioral time scales: experimental support of neohebbian three-factor learning rules. *Frontiers in neural circuits*, 12:53, 2018.
- Shixiang Gu, Ethan Holly, Timothy Lillicrap, and Sergey Levine. Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates. In 2017 IEEE international conference on robotics and automation (ICRA), pp. 3389–3396. IEEE, 2017.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–1870. Pmlr, 2018.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456. pmlr, 2015.
- Haiyan Jiang, Vincent Zoonekynd, Giulia De Masi, Bin Gu, and Huan Xiong. Tab: Temporal accumulated batch normalization in spiking neural networks. In *The Twelfth International Conference on Learning Representations*, 2024.
- Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. 1960.
- Youngeun Kim and Priyadarshini Panda. Revisiting batch normalization for training low-latency deep spiking neural networks from scratch. *Frontiers in neuroscience*, 15:773954, 2021.
- Jens Kober, J Andrew Bagnell, and Jan Peters. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32(11):1238–1274, 2013.
- Vijay Konda and John Tsitsiklis. Actor-critic algorithms. *Advances in neural information processing systems*, 12, 1999.
- Aakash Kumar, Lei Zhang, Hazrat Bilal, Shifeng Wang, Ali Muhammad Shaikh, Lu Bo, Avinash Rohra, and Alisha Khalid. Dsqn: Robust path planning of mobile robot based on deep spiking q-network. *Neurocomputing*, 634:129916, 2025.
- Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv* preprint arXiv:1509.02971, 2015.
- TP Lillicrap. Continuous control with deep reinforcement learning. arXiv preprint arXiv:1509.02971, 2015.
- Guisong Liu, Wenjie Deng, Xiurui Xie, Li Huang, and Huajin Tang. Human-level control through directly trained deep spiking q-networks. *IEEE transactions on cybernetics*, 53(11):7187–7198, 2022.
- Wolfgang Maass. Networks of spiking neurons: the third generation of neural network models. *Neural networks*, 10(9):1659–1671, 1997.
- Volodymyr Mnih. Playing atari with deep reinforcement learning. arXiv preprint arXiv:1312.5602, 2013.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015a.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015b.
- Devdhar Patel, Hananel Hazan, Daniel J Saunders, Hava T Siegelmann, and Robert Kozma. Improved robustness of reinforcement learning policies upon conversion to spiking neuronal network platforms applied to atari breakout game. *Neural Networks*, 120:108–115, 2019.

- Lang Qin, Rui Yan, and Huajin Tang. A low latency adaptive coding spiking framework for deep reinforcement learning. *arXiv* preprint arXiv:2211.11760, 2022.
- Shibani Santurkar, Dimitris Tsipras, Andrew Ilyas, and Aleksander Madry. How does batch normalization help optimization? *Advances in neural information processing systems*, 31, 2018.
- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv* preprint arXiv:1506.02438, 2015.
- Yinqian Sun, Yi Zeng, and Yang Li. Solving the spike feature information vanishing problem in spiking deep q network with potential based normalization. *Frontiers in Neuroscience*, 16:953368, 2022.
- Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction. MIT press, 2018.
- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.
- Weihao Tan, Devdhar Patel, and Robert Kozma. Strategy and benchmark for converting deep q-networks to event-driven spiking neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pp. 9816–9824, 2021.
- Guangzhi Tang, Neelesh Kumar, and Konstantinos P Michmizos. Reinforcement co-learning of deep and spiking neural networks for energy-efficient mapless navigation with neuromorphic hardware. In 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 6090– 6097. IEEE, 2020.
- Guangzhi Tang, Neelesh Kumar, Raymond Yoo, and Konstantinos Michmizos. Deep reinforcement learning with population-coded spiking neural network for continuous control. In *Conference on Robot Learning*, pp. 2016–2029. PMLR, 2021.
- Emanuel Todorov. Convex and analytically-invertible dynamics with contacts and constraints: Theory and implementation in mujoco. In 2014 IEEE International Conference on Robotics and Automation (ICRA), pp. 6054–6061. IEEE, 2014a.
- Emanuel Todorov. Convex and analytically-invertible dynamics with contacts and constraints: Theory and implementation in mujoco. In 2014 IEEE International Conference on Robotics and Automation (ICRA), pp. 6054–6061. IEEE, 2014b.
- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In 2012 IEEE/RSJ international conference on intelligent robots and systems, pp. 5026–5033. IEEE, 2012.
- Mark Towers, Ariel Kwiatkowski, Jordan K Terry, John U Balis, Gianluca De Cola, Tristan Deleu, Manuel Goulão, Andreas Kallinteris, Markus Krimmel, KG Arjun, et al. Gymnasium: A standard interface for reinforcement learning environments. *CoRR*, 2024.
- Paweł Wawrzyński. A cat-like robot real-time learning to run. In Adaptive and Natural Computing Algorithms: 9th International Conference, ICANNGA 2009, Kuopio, Finland, April 23-25, 2009, Revised Selected Papers 9, pp. 380–390. Springer, 2009.
- Yujie Wu, Lei Deng, Guoqi Li, Jun Zhu, and Luping Shi. Spatio-temporal backpropagation for training high-performance spiking neural networks. *Frontiers in neuroscience*, 12:331, 2018.
- Yuxin Wu and Justin Johnson. Rethinking" batch" in batchnorm. arXiv preprint arXiv:2105.07576, 2021.
- Zijie Xu, Tong Bu, Zecheng Hao, Jianhao Ding, and Zhaofei Yu. Proxy target: Bridging the gap between discrete spiking neural networks and continuous control. *arXiv preprint arXiv:2505.24161*, 2025.

- Kashu Yamazaki, Viet-Khoa Vo-Ho, Darshan Bulsara, and Ngan Le. Spiking neural networks and their applications: A review. *Brain sciences*, 12(7):863, 2022.
- Zhile Yang, Shangqi Guo, Ying Fang, Zhaofei Yu, and Jian K Liu. Spiking variational policy gradient for brain inspired reinforcement learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- Duzhen Zhang, Tielin Zhang, Shuncheng Jia, and Bo Xu. Multi-sacle dynamic coding improved spiking actor network for reinforcement learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pp. 59–67, 2022.
- Duzhen Zhang, Qingyu Wang, Tielin Zhang, and Bo Xu. Biologically-plausible topology improved spiking actor network for efficient deep reinforcement learning. *arXiv preprint arXiv:2403.20163*, 2024.
- Hanle Zheng, Yujie Wu, Lei Deng, Yifan Hu, and Guoqi Li. Going deeper with directly-trained larger spiking neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pp. 11062–11070, 2021.