# GRSN: Gated Recurrent Spiking Neurons for POMDPs and MARL

**Lang Qin**[1] , **Ziming Wang**[1] , **Runhao Jiang**[1] , **Rui Yan**[2] and **Huajin Tang**[1*]

[1]College of Computer Science and Technology, Zhejiang University, Hangzhou, China
[2]College of Computer Science, Zhejiang University of Technology, Hangzhou, China

qinl@zju.edu.cn, zi_ming_wang@outlook.com,
15520816169@163.com, ryan@zjut.edu.cn, htang@zju.edu.cn

## Abstract

Spiking neural networks (SNNs) are widely applied in various fields due to their energy-efficient and fast-inference capabilities. Applying SNNs to reinforcement learning (RL) can significantly reduce the computational resource requirements for agents and improve the algorithm's performance under resource-constrained conditions. However, in current spiking reinforcement learning (SRL) algorithms, the simulation results of multiple time steps can only correspond to a single-step decision in RL. This is quite different from the real temporal dynamics in the brain and also fails to fully exploit the capacity of SNNs to process temporal data. In order to address this temporal mismatch issue and further take advantage of the inherent temporal dynamics of spiking neurons, we propose a novel temporal alignment paradigm (TAP) that leverages the single-step update of spiking neurons to accumulate historical state information in RL and introduces gated units to enhance the memory capacity of spiking neurons. Experimental results show that our method can solve partially observable Markov decision processes (POMDPs) and multi-agent cooperation problems with similar performance as recurrent neural networks (RNNs) but with about 50% power consumption.

## 1 Introduction

Spiking neural networks (SNNs) possess unique spatiotemporal dynamics and all-or-none firing characteristics, enabling them to perform low-power and high-speed inference on neuromorphic hardware [Roy *et al.*, 2019; Pei *et al.*, 2019]. With the continuous improvement of SNN training algorithms, the learning capability and network scalability of SNNs are gradually approaching those of artificial neural networks (ANNs) [Lian *et al.*, 2023; Wang *et al.*, 2023; Qin *et al.*, 2023a]. Therefore, in recent years, there have been more and more efforts to explore SNNs as alternatives to ANNs in the field of deep reinforcement learning (DRL),
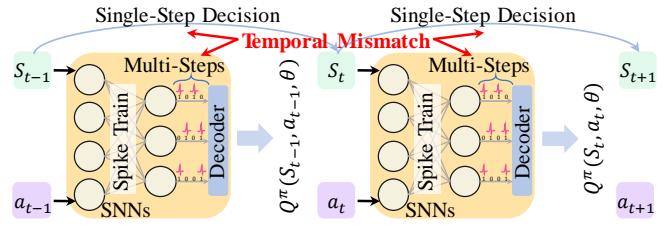
---

*[*]Corresponding author



Figure 1: Temporal mismatch issue in SRL algorithms. The figure shows the basic model of SRL methods; it uses spike results of multiple time steps to decode for one-step value function calculation or action selection in RL.

aiming to reduce the computational cost and inference latency of intelligent agents.

In SNN learning algorithms, firing rates are often used to calculate continuous values due to the discrete nature of spikes. However, the firing rates, which are in the range of 0 to 1, are difficult to map onto continuous values (such as continuous actions and Q-values) in RL that do not have value range limitations. There are generally three types of approaches to address or avoid this issue in order to achieve deep spiking reinforcement learning (SRL): (1) ANN-SNN conversion methods for RL; (2) a hybrid framework for the actor-critic (AC) method; and (3) directly trained SNNs for RL.

The conversion algorithms [Patel *et al.*, 2019; Tan *et al.*, 2021] avoid the aforementioned mapping problem by using a pre-trained ANN-RL model to transform the SNN model, thus avoiding direct involvement in the learning process of RL algorithms. However, the conversion algorithms often require high simulation time steps due to the computational precision requirement, resulting in increased energy consumption and inference delay. To reduce simulation steps, a hybrid framework [Tang *et al.*, 2020a; Tang *et al.*, 2020b; Zhang *et al.*, 2022] for the AC method has been proposed. This approach combines spiking actor networks and artificial critic networks and utilizes a gradient descent algorithm for collaborative training. The critic networks that require precise calculations are completed by ANN, while SNN is only used for selecting actions. This hybrid framework successfully reduces the need for simulation time steps, thereby reducing energy consumption and inference delay.

Clearly, the hybrid framework method is only applicable to RL algorithms based on the AC method. Additionally, similar to conversion algorithms, the hybrid framework method still relies on ANN to accomplish RL tasks. To enhance the generalizability of the algorithms and eliminate reliance on ANN, directly trained SNNs for RL [Chen *et al.*, 2022; Liu *et al.*, 2022; Sun *et al.*, 2022; Qin *et al.*, 2023b] have been proposed. These algorithms directly train SNNs through surrogate gradient learning [Neftci *et al.*, 2019], and use suitable coders to handle computations involving continuous values. These methods do not rely on ANNs and can complete tasks in a few time steps. However, such methods rely more on the design and selection of coders and also have a trade-off between performance and latency.

Reinforcement learning algorithms are fundamentally designed to address sequential decision problems; they can synergize effectively with SNNs that inherently capture temporal sequence dynamics. However, the aforementioned existing SRL methods still use multiple time steps to decode the spike information to obtain continuous values, which leads to a serious temporal mismatch problem (Figure 1). This temporal mismatch between multiple time steps and single-step decisions not only deviates from the real dynamics in the brain but also impairs the ability of SNNs to process temporal data.

To address this issue, we propose a novel SRL paradigm that aligns the single-step updating of spiking neurons with the single-step decisions in RL, enabling the sequence decision problem to be solved within a whole simulated time window. Considering the weak temporal correlations between time steps of the original spiking neurons and the inability to guarantee accuracy after temporal alignment, we further introduced gated units to enhance the long- and short- term memory capabilities of neurons, ensuring the effectiveness of the proposed paradigm. In order to test the performance of the proposed gated recurrent spiking neurons (GRSN) and temporal alignment paradigm (TAP), we conducted experiments in partially observable (PO) and multi-agent environments. Our main contributions are summarized as follows:

- We pointed out the issue of temporal mismatch in current SRL algorithms and proposed a novel temporal alignment paradigm (TAP) for SRL to solve this issue.

- We designed a novel gated recurrent spiking neuron (GRSN) with enhanced long- and short-term memory capabilities, which can be applied to the proposed TAP.

- To our best knowledge, this work is the first to use SNNs to address POMDPs and multi-agent reinforcement learning (MARL) problems. We verified the effectiveness of SNNs in PO and multi-agent environments.

- Experimental results show that GRSN can outperform original spiking neurons in benchmark environments and can achieve similar performance as RNN-RL with about 50% energy consumption.

## 2 Related Works

### 2.1 Spiking Reinforcement Learning

The development of SRL can be mainly divided into three periods: (1) the basic research of SNNs and RL [Williams, 1992; Bohté *et al.*, 2002]; (2) the synaptic-plasticity-based SRL algorithm [Seung, 2003; Urbanczik and Senn, 2009; Frémaux *et al.*, 2010]; and (3) the combination of deep RL and SNNs [Zhang *et al.*, 2021; Liu *et al.*, 2022; Qin *et al.*, 2023b]. Early-stage works often study the combination of synaptic plasticity and RL theory, which aims to reveal how reward mechanisms in the brain correlate and combine with RL algorithms. In the later stage, with the development of DRL, the SRL algorithm focuses on applying SNNs to DRL, aiming for better performance.

### 2.2 Recurrent Spiking Neural Networks

Due to the inherent temporal nature of spiking neurons, there have been many attempts to use recurrent spiking neural networks for temporal data processing. Some work [Lotfi-Rezaabad and Vishwanath, 2020; Rao *et al.*, 2022] has designed spike-based long short-term memory networks and tested them on audio or text datasets. Other work [Yin *et al.*, 2020; Ponghiran and Roy, 2022; Liu *et al.*, 2023] is to explore the long-term dependence of spiking neurons and complete structure optimization or time series forecasting tasks.

## 3 Preliminaries

### 3.1 Notations

For all variables, superscripts refer to the number of layers in the neural network, and subscripts refer to their time steps in the time window. For example, $u_t^l$ denotes the membrane potential in layer $l$ at the $t$-th time step. We follow the conventions representing vectors and matrix with bold italic letters and bold capital letters, respectively, such as $o$ and $W$. We use $\odot$ to signify element-wise multiplication for two vectors.

### 3.2 Leaky Integrate-and-Fire Model

Different from traditional neural networks, SNNs use binary spikes as information carriers. In order to compensate for the lack of binary spikes in information expression, the time dimension, or latency, is introduced into SNNs. SNNs accept event streams as input, and the forward propagation of SNNs is repeated for $T$ time steps in the temporal dimension to calculate the output spike trains. Therefore, the basic computing unit of SNNs, spiking neurons, has unique spatio-temporal dynamics. Here we introduce the widely used Leaky Integrate-and-Fire (LIF) neuron [Wu *et al.*, 2018] as the benchmark in this work. In each time step $t$, the membrane potential $u_t^l$ of spiking neurons at the $l$-th layer will integrate the input current $c_t^l$ and the decay voltage $\beta \hat{u}_{t-1}^l$:

$$u_t^l = \beta \hat{u}_{t-1}^l + (1 - \beta) c_t^l \tag{1}$$

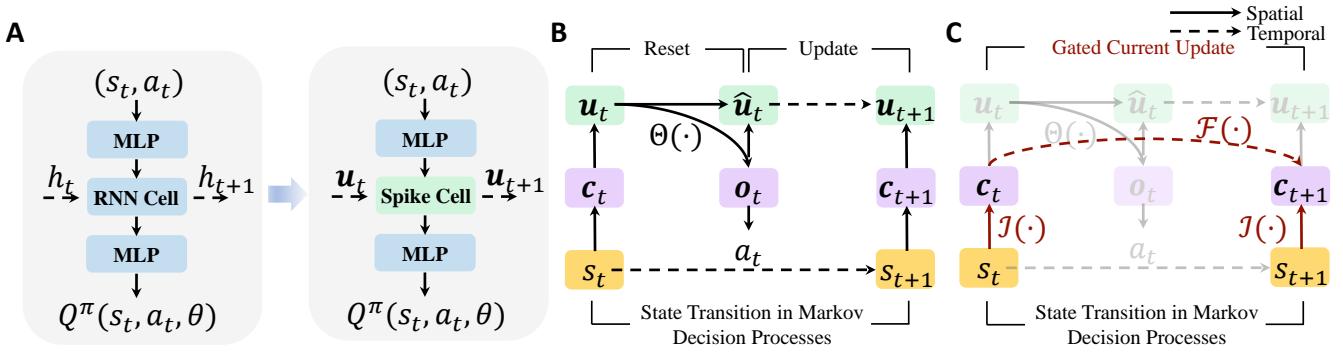$$c_t^l = W^l o_t^{l-1} + b^l \tag{2}$$

Figure 2: Overview of the proposed TAP and GRSN. **A.** Basic model of RNN-RL and SNN-RL algorithms under temporal alignment paradigm. Spike cell is a special form of single-step update for spiking neurons, which maintains the update of membrane potential $u_t$ and outputs the $Q$ function. The multilayer perceptron (MLP) serves as the embedding or coding layer here. **B.** Detailed update dynamics for the parameters of spiking neurons. The state $s_t$ is encoded into the input current $c_t$ and feedforward to the spiking neurons. The neurons calculate the current membrane potential $u_t$ based on $c_t$ and $u_{t-1}$, and determine whether to emit spikes $o_t$. The output spikes are decoded and used to obtain the action $a_t$, and then the next state $s_{t+1}$ is achieved based on the state transition function. **C.** Dynamics of GRSN. The red part shows the added recurrent connection and gated unit. $\mathcal{F}$ and $\mathcal{I}$ represent the forget gates and input gates, respectively. GRSN also follows the temporal alignment paradigm, with each time step corresponding to a single-step state transition.

Membrane potential will be reset to $u_r$ (default set to 0) when a spike is emitted at the last time step:

$$\hat{u}_t^l = u_r o_t^l + (1 - o_t^l) u_t^l \qquad (3)$$

The neurons will emit output spikes $o_t^l$ whenever the membrane potential $u_t^l$ exceeds the threshold $\vartheta$:

$$o_t^l = \Theta(u_t^l - \vartheta) \qquad (4)$$

where $\Theta(x)$ is the Heaviside function:

$$\Theta(x) = \begin{cases} 1, if \quad x \geq 0 \\ 0, otherwise \end{cases} \qquad (5)$$

### 3.3 Surrogate Gradient

Due to the non-differentiability of the binary firing function $\Theta(x)$, SNNs often employ surrogate gradients [Neftci *et al.*, 2019] during back-propagation (BP). We adopt the arctangent function to replace the derivative of the binary firing function:

$$\Theta'(x) \triangleq h'(x) = \frac{\alpha}{2[1 + (\frac{\pi}{2}\alpha x)^2]} \qquad (6)$$

where $\alpha$ is a hyper-parameter that controls the width of the surrogate function. With the surrogate gradients for binary firing functions, we can directly use the BP method to train SNNs.

## 4 Method

### 4.1 Temporal Alignment Paradigm

Reinforcement learning tasks can be modeled as Markov decision processes (MDP) $(\mathcal{S}, \mathcal{A}, \mathcal{R}, p, \gamma)$, with state space $\mathcal{S}$, action space $\mathcal{A}$, scalar reward function $\mathcal{R}$, transition dynamics $p$, and discount factor $\gamma$ [Sutton and Barto, 2018]. RL agents are controlled by policy $\pi$, with the goal of maximizing the expected discounted return $\mathbb{E}_\pi[\sum_{t=0}^{\infty} \gamma^t r_{t+1}]$. According to

the Bellman equation [Bellman, 1966], we can obtain the iteration form of the value function:

$$V^\pi(s_t) = \mathbb{E}_\pi[R_{t+1} + \gamma V^\pi(s_{t+1})] \qquad (7)$$

When the policy is determined, Eq 7 can be simplified as:

$$V(s_t) = \frac{1}{\gamma} V(s_{t-1}) + \frac{-1}{\gamma} \mathbb{E}[R_t] \qquad (8)$$

Observing Eq 8 and Eq 1, we can find that the state transitions in MDPs exhibit a striking resemblance to the state changes of spiking neurons. Therefore, the key to solving the temporal mismatch problem is to utilize this similarity.

We propose a novel SRL paradigm according to this similarity: within the framework of RL algorithms, the application of SNNs in a manner similar to RNNs (Figure 2A), aligning each step of spiking neuron updates $u_t \rightarrow u_{t+1}$ with each state transition $s_t \rightarrow s_{t+1}$ in an MDP. The advantage of this approach lies in the natural utilization of the temporal structure of SNNs while significantly reducing the number of time steps required for the entire task. Figure 2A and 2B show the proposed TAP and the detailed dynamics of spiking neurons.

### 4.2 Gradient Analysis of Spiking Neurons

In Section 3.2, we introduce the basic computational unit of SNNs: LIF neurons. It uses a constant leaky factor $\beta$ to play the roles of forget gate and input gate. The use of this constant factor results in weaker temporal correlations in the neuron's internal state ($u_t$), making it challenging to handle data with long-term dependencies. In order to enhance the memory capacity of SNNs in the temporal domain, we need to analyze and optimize the LIF neuron.

According to Eq 1-5, it can be observed that the LIF neuron has two internal states: membrane potential $u_t^l$ and input current $c_t^l$. To compare the importance of these two internal states in backpropagation, we need to calculate the gradient

of the loss function for connection weight $\boldsymbol{W}^l$. Suppose $\mathcal{L}$ is the loss function that we would like to minimize. According to Eq 1-5 and the chain rule, its gradient is:

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{W}^l} = \sum_{t=1}^{T} \frac{\partial \mathcal{L}}{\partial \boldsymbol{o}_t^l} \frac{\partial \boldsymbol{o}_t^l}{\partial \boldsymbol{u}_t^l} \frac{\partial \boldsymbol{u}_t^l}{\partial \boldsymbol{W}^l} \tag{9}$$

$T$ is the total number of time steps. The first term $\frac{\partial \mathcal{L}}{\partial \boldsymbol{o}_t^l}$ in Eq 9 is determined by the decoder and loss function and does not affect the flow of the gradient. According to Eq 4 and Eq 6, the second term could be derived as:

$$\frac{\partial \boldsymbol{o}_t^l}{\partial \boldsymbol{u}_t^l} = \frac{\partial \Theta(\boldsymbol{u}_t^l - \vartheta)}{\partial \boldsymbol{u}_t^l} = h'(\boldsymbol{u}_t^l - \vartheta) \tag{10}$$

Then we calculate the third term $\frac{\partial \boldsymbol{u}_t^l}{\partial \boldsymbol{W}^l}$:

$$\frac{\partial \boldsymbol{u}_t^l}{\partial \boldsymbol{W}^l} = \frac{\partial \boldsymbol{u}_t^l}{\partial \boldsymbol{c}_t^l} \frac{\partial \boldsymbol{c}_t^l}{\partial \boldsymbol{W}^l} + \frac{\partial \boldsymbol{u}_t^l}{\partial \boldsymbol{u}_{t-1}^l} \frac{\partial \boldsymbol{u}_{t-1}^l}{\partial \boldsymbol{W}^l} \tag{11}$$

Bring Eq 1 into Eq 11:

$$\frac{\partial \boldsymbol{u}_t^l}{\partial \boldsymbol{W}^l} = [(1-\beta) + \beta \frac{\partial \hat{\boldsymbol{u}}_{t-1}^l}{\partial \boldsymbol{c}_t^l}] \frac{\partial \boldsymbol{c}_t^l}{\partial \boldsymbol{W}^l} + \\ \beta \frac{\partial [u_r \boldsymbol{o}_{t-1}^l + (1 - \boldsymbol{o}_{t-1}^l)\boldsymbol{u}_{t-1}^l]}{\partial \boldsymbol{u}_{t-1}^l} \frac{\partial \boldsymbol{u}_{t-1}^l}{\partial \boldsymbol{W}^l} \tag{12}$$

Obviously, $\frac{\partial \hat{\boldsymbol{u}}_{t-1}^l}{\partial \boldsymbol{c}_t^l} = 0$. Let

$$\delta_t = \frac{\partial \hat{\boldsymbol{u}}_t^l}{\partial \boldsymbol{u}_t^l} = \frac{\partial [u_r \boldsymbol{o}_t^l + (1 - \boldsymbol{o}_t^l)\boldsymbol{u}_t^l]}{\partial \boldsymbol{u}_t^l} \\ = [(u_r - \boldsymbol{u}_t^l)h'(\boldsymbol{u}_t^l - \vartheta) + 1 - \boldsymbol{o}_t^l] \tag{13}$$

Then Eq 12 can be simplified as:

$$\frac{\partial \boldsymbol{u}_t^l}{\partial \boldsymbol{W}^l} = (1-\beta)\frac{\partial \boldsymbol{c}_t^l}{\partial \boldsymbol{W}^l} + \beta \delta_{t-1} \frac{\partial \boldsymbol{u}_{t-1}^l}{\partial \boldsymbol{W}^l} \tag{14}$$

Continuing substitution $\frac{\partial \boldsymbol{u}_t^l}{\partial \boldsymbol{W}^l}$ leads to a closed-form expression:

$$\frac{\partial \boldsymbol{u}_t^l}{\partial \boldsymbol{W}^l} = (1-\beta)\sum_{i=1}^{t} \beta^i \prod_{j=1}^{i} \delta_{t-j} \frac{\partial \boldsymbol{c}_{t-i}^l}{\partial \boldsymbol{W}^l} + \\ \prod_{k=1}^{t} \beta \delta_{t-k} + (1-\beta)\frac{\partial \boldsymbol{c}_t^l}{\partial \boldsymbol{W}^l} \tag{15}$$

When the variable subscript $t = 0$, its value is the artificially set initial value.

To analyze the main carriers of gradient, we propose the following theorem:

**Theorem 1.** *With the LIF model (Eq 1-5) and the common parameter settings ($\beta = 0.5, \vartheta = 1, \alpha = 2, u_r = 0$), the gradient of loss function $\frac{\partial \mathcal{L}}{\partial \boldsymbol{W}^l}$ mainly depends on current $\frac{\partial \boldsymbol{c}}{\partial \boldsymbol{W}^l}$ when total time-step $T$ is large enough.*

*Proof.* The first term in Eq 9 depends on the decoding method, and the second term in Eq 9 depends on the SG
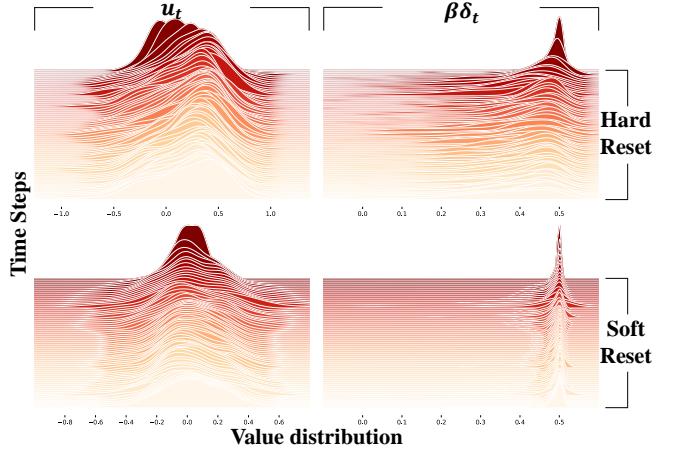


Figure 3: Distributions of membrane potential $\boldsymbol{u}_t$ and the term $\beta \delta_t$. The y-axis represents the time step, and the x-axis represents the value of the parameter. The probability density curve corresponding to each time step represents the parameter distribution of all neurons in the single-layer network.

function. Both of them do not affect the direction of gradient flow. Therefore, the third term ($\frac{\partial \boldsymbol{u}_t^l}{\partial \boldsymbol{W}^l}$) is the most significant. By analyzing the monotonicity (presented in the appendix) of Eq 15, it can be found that when $T$ is large, $\prod_{i=1}^{T} \beta \delta_t$ tends to 0. Hence, the last term $(1-\beta)\frac{\partial \boldsymbol{c}_t^l}{\partial \boldsymbol{W}^l}$ in Eq 15 determines the direction of the gradient. $\square$

To further verify **Theorem** 1, we conducted experiments in *CartPole-V* (see Section 5) environment and visualized the parameter distribution during the training process (Figure 3). From figure 3, the maximum absolute value of the term $\beta \delta_t$ is less than 1, and most of the membrane potential values are distributed around 0. Hence, $\prod_{i=1}^{T} \beta_t$ will tend to 0 when $T$ is large enough. This result is consistent with **Theorem** 1.

According to **Theorem** 1, the gradient of the loss function $\mathcal{L}$ is mainly transmitted through the input current $\boldsymbol{c}$. Therefore, adding gated recurrent connections to the input current $\boldsymbol{c}$ can enhance the memory capability of spiking neurons more effectively than membrane potential $\boldsymbol{u}$.

### 4.3 Gated Recurrent Spiking Neurons

Optimizing the dynamics of neurons and improving memory ability is mainly divided into three steps. The first step is to change the decay factor $\beta$ to a learnable parameter, which can be jointly optimized with the entire network. The learnable decay factor can dynamically adjust the forgetting ability of neurons during the training process and adaptively find the optimal value [Rathi and Roy, 2023]. Then, we use a soft reset instead of a hard reset (**Theorem** 1 still holds in the soft reset setting; see Figure 3 and appendix for details). Compared to the hard reset method, soft reset can reduce temporal information loss, and it has better performance when the network is shallow [Ledinauskas *et al.*, 2020]. Finally, we added recurrent connections and adopted gating functions to enhance the temporal correlation of the neuron's internal states $\boldsymbol{c}$.

| Env. Net. | Learning Paradigm | CartPole | | Pendulum | |
|---|---|---|---|---|---|
| | | -V | -P | -V | -P |
| MLP [*] | Actor-Critic | $21.6 \pm 1.49$ | $21.2 \pm 1.95$ | $-1502.8 \pm 61.09$ | $-1380.1 \pm 47.45$ |
| GRU [*] | RNN-RL | $\mathbf{200.0 \pm 0.00}$ | $196.7 \pm 6.30$ | $\mathbf{-181.3 \pm 15.47}$ | $-203.5 \pm 16.36$ |
| LIF | **TAP** | $\mathbf{200.0 \pm 0.00}$ | $152.2 \pm 21.21$ | $-942.4 \pm 92.19$ | $-824.8 \pm 178.80$ |
| **GRSN** | **TAP** | $\mathbf{200.0 \pm 0.00}$ | $199.9 \pm 0.20$ | $-189.4 \pm 17.64$ | $\mathbf{-195.8 \pm 43.99}$ |

[*] Reproduction results of [Ni *et al.*, 2022].

Table 1: Performance Comparison on PO Environment

---

**Algorithm 1** GRSN with temporal alignment paradigm

1: **Input:** States sequence $< s_1, s_2, ..., s_T >$
2: **Parameter:** Total steps $T$, total layers $L$
3: **Output:** Values sequence $< q_1, q_2, ..., q_T >$
4: initialize the hidden states $\boldsymbol{u}_0^1 = \hat{\boldsymbol{u}}_0^1 = \mathbf{0}, \boldsymbol{c}_0^1 = \mathbf{0}$
5: **for** $t = 1$ to $T$ **do**
6:     initialize the input spike $\boldsymbol{o}_t^0 = encoder(s_t)$
7:     **for** $l = 1$ to $L$ **do**
8:         update $\boldsymbol{c}_t^l$ and $\boldsymbol{u}_t^l$ // Eq 17, Eq 1
9:         firing and computing output $\boldsymbol{o}_t^l$ // Eq 4
10:       reset the potential $\hat{\boldsymbol{u}}_t^l$ // Eq 16
11:     **end for**
12:     compute the values $q_t = decoder(\boldsymbol{o}_t^L)$
13: **end for**
14: **return** $< q_1, q_2, ..., q_T >$

**Notes**: The encoder and decoder are implemented by MLP, similar to the embedding layers in RL.

---

We named the proposed neuron model gated recurrent spiking neurons (GRSN). Its membrane potential update rules, firing rules, and definitions of Heaviside function are consistent with those of LIF neurons (see Eq 1, Eq 4 and Eq 5). Different from LIF neurons (Eq. 3), the reset rule of GRSN is soft reset:

$$\hat{\boldsymbol{u}}_t^l = \boldsymbol{u}_t^l - \vartheta \boldsymbol{o}_t^l \tag{16}$$

GRSN also adds recurrent connections at the current level:

$$\boldsymbol{c}_t^l = \mathcal{F}(\boldsymbol{o}_t^{l-1}) \odot \boldsymbol{c}_{t-1}^l + [1 - \mathcal{F}(\boldsymbol{o}_t^{l-1})] \odot \mathcal{I}(\boldsymbol{o}_t^{l-1}) \tag{17}$$

where $\mathcal{F}(\cdot)$ is the forgetting gate that regulates the proportion of forget and input, and $\mathcal{I}(\cdot)$ is the input gate that processes input signals. Both $\mathcal{F}(\cdot)$ and $\mathcal{I}(\cdot)$ are composed of a fully connected (FC) layer and an activation function. Their detailed definitions are shown as follows:

$$\begin{aligned} \mathcal{F}(\boldsymbol{x}) &= \sigma(\boldsymbol{W}_f \boldsymbol{x} + \boldsymbol{b}_f) \\ \mathcal{I}(\boldsymbol{x}) &= ReLU(\boldsymbol{W}_i \boldsymbol{x} + \boldsymbol{b}_i) \end{aligned} \tag{18}$$

where $\sigma$ is the sigmoid function.

Figure 2C shows the dynamics of GRSN in the temporal alignment paradigm. Compared with Figure 2B, red lines indicate the recurrent connections. The pseudocode of GRSN with the temporal alignment paradigm is shown in Algorithm 1.

## 5 Experiments

In the TAP, SNN behaves more like an RNN rather than a traditional DNN. Spiking neurons calculate the value function by accumulating information from multiple previous consecutive states. This configuration is often employed in RL to address POMDPs. Consequently, to validate the effectiveness and versatility of the proposed GRSN, we conducted experiments in both PO environments and multi-agent cooperative environments. The detailed parameters and settings of the experiment can be found in the appendix.

### 5.1 Partially Observable Classic Control Tasks

**Enviroments & Experimental Settings**

The *Pendulum* and *CartPole* tasks are classic control tasks for evaluating RL algorithms. The cart-pole problem described in [Barto *et al.*, 1983] aims to balance the pole by applying forces in the left and right directions on the cart (Figure 4A). The inverted pendulum swingup problem is based on the classic problem in control theory. Its goal is to apply torque on the free end to swing it into an upright position, with its center of gravity right above the fixed point (Figure 4A). We conducted experiments in PO cases of these two control tasks [Han *et al.*, 2020], in which only velocities or positions could be observed. We use suffix names (-P/-V) to distinguish these two settings.

We follow the experimental settings proposed in [Ni *et al.*, 2022] and use SAC [Haarnoja *et al.*, 2018] for discrete environments and TD3 [Fujimoto *et al.*, 2018] for continuous environments. For networks, we use single-layer gated recurrent units (GRU) [Chung *et al.*, 2014] for RNNs, two fully connected layers for MLP, and single-layer LIF/GRSN for SNNs. We conducted independent experiments on five different random seeds for each method and tested each final model ten times to eliminate the interference of randomness. Experiment details can be found in the appendix.

**Results**

To demonstrate the superiority of the TAP and GRSN in multiple aspects, we selected three different networks and paradigm settings as baselines (Table 1). From the results, TAP effectively utilizes temporal information, and its performance far exceeds that of traditional algorithms using MLP. On the other hand, GRSN enhances the temporal correlation of spiking neurons, thereby further improving performance. The combination of TAP and GRSN equals or even surpasses the state-of-the-art RNN-RL method. From the learning curve (Figure 5), it can be seen that the convergence rate

| Senarios | Difficulty | QMIX-GRU [†] | | QMIX-GRSN | |
|---|---|---|---|---|---|
| | | Mean reward | Win rate (%) | Mean reward | Win rate (%) |
| *8m* | Easy | $19.8 \pm 0.26$ | $97.6 \pm 3.49$ | **$20.0 \pm 0.12$** | **$99.4 \pm 1.48$** |
| *2s3z* | Easy | **$19.9 \pm 0.18$** | **$97.9 \pm 2.84$** | $19.9 \pm 0.20$ | $97.0 \pm 3.57$ |
| *8m_vs_9m* | Hard | $19.4 \pm 0.69$ | $91.6 \pm 8.73$ | **$19.6 \pm 0.35$** | **$94.6 \pm 4.54$** |
| *3s_vs_5z* | Hard | $21.0 \pm 0.44$ | **$97.1 \pm 3.46$** | **$21.4 \pm 0.60$** | $93.0 \pm 5.48$ |
| *27m_vs_30m* | Super Hard | $19.2 \pm 0.51$ | $79.4 \pm 12.19$ | **$19.9 \pm 0.14$** | **$96.4 \pm 3.74$** |
| *MMM2* | Super Hard | $17.6 \pm 3.76$ | $75.2 \pm 37.93$ | **$18.5 \pm 0.73$** | **$83.6 \pm 8.17$** |

[†] Reproduction results of [Hu *et al.*, 2021].

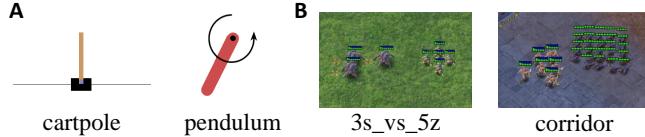Table 2: Experimental results in SMAC environment

Figure 4: **A.** Diagram of the pendulum and cartpole problem. The agent needs to apply appropriate forces to control the system to reach the target state, in order to get rewards. **B.** Game scenarios of two different SMAC maps. Agents need to defeat enemies to achieve high rewards and win the game.

of GRSN is slightly weaker than that of traditional RNNs in some environments due to the use of discrete spikes as information carriers, but in the end, the performance of the two is very similar. LIF and MLP, which are baselines, are weaker than GRSN and RNNs in terms of temporal correlation, so their performance and convergence rate are far inferior to the other two networks.

## 5.2 StarCraft Multi-Agent Challenge

### Enviroments & Experimental Settings

StarCraft Multi-Agent Challenge (SMAC) is a benchmark environment for evaluating MARL algorithms [Samvelyan *et al.*, 2019]. SMAC is based on the popular real-time strategy (RTS) game StarCraft II, and it contains multiple micro StarCraft II scenarios (Figure 4B). In SMAC, the overall goal is to maximize the win rate for each battle scenario. Each agent in SMAC has a circular field of view centered on itself. With such a limited field of view, each agent is faced with a partially observable environment. The action space is discrete and includes *move[direction]*, *attack[enemy_id]*, *stop*, and *no-op*. Healing units use *heal[agent_id]* actions instead of *attack[enemy_id]*. The SMAC reward includes several parts: hit-point damage dealt and enemy units killed, together with a special bonus for winning the battle.

The QMIX [Rashid *et al.*, 2018] is a classic value-based multi-agent collaboration algorithm. QMIX adopts the centralized training distributed execution (CTDE) mode, and each agent uses RNNs as main networks. Therefore, the proposed GRSN and TAP can be well adapted to the QMIX algorithm. We follow the experimental and parameter settings of QMIX and select SMAC maps of different difficulty levels for the experiment. Detailed parameters and settings are also listed in the appendix.

## Results

To evaluate the performance of the method more comprehensively, we selected several SMAC maps with different difficulty levels for experiments. We conducted five independent experiments in each environment, training 10 million steps per experiment. We use the GRU-based QMIX algorithm as the baseline, and the experimental results are shown in Table 2. From the experimental results, it can be seen that the GRSN-based QMIX algorithm surpasses the original GRU-based QMIX algorithm in most environments, which proves that GRSN achieves or even exceeds GRU in the processing of temporal information. In addition, it is not difficult to find that the advantages of GRSN are more obvious in difficult scenarios. This also proves that GRSNs using spikes have better robustness than traditional RNNs.

## 5.3 Ablation Study

In Section 4, we propose the TAP that can solve the temporal mismatch problem and greatly reduce time steps, and the GRSN that can enhance the temporal association of spiking neurons and improve performance. To demonstrate the impact of these two on the training speed and algorithm performance more intuitively, we conducted multiple groups of ablation experiments. We conducted experiments in both classical control tasks and SMAC environments. Similar to the previous setting, we conducted five independent experiments at each of the different settings and averaged them as the final results (Table 3).

### Temporal Alignment Paradigm

The TAP solves the mismatch problem, thus significantly reducing the number of time steps and training and inference time. From Table 3, the algorithm using TAP reduces the time step by $T$ times, where $T = 4$ is the time window size of the basic spiking neuron. In addition, TAP has also shortened the training time by about half while maintaining performance, greatly improving the efficiency of algorithms.

### GRSN vs. LIF

In Table 3, we can compare the GRSN and LIF neurons. GRSN has better performance in different environmental and paradigm settings. However, due to the additional gating function, the training time of GRSN will be slightly longer than that of LIF. In terms of stability, GRSN also demonstrated the most stable performance, with the lowest standard deviation in five independent experiments.
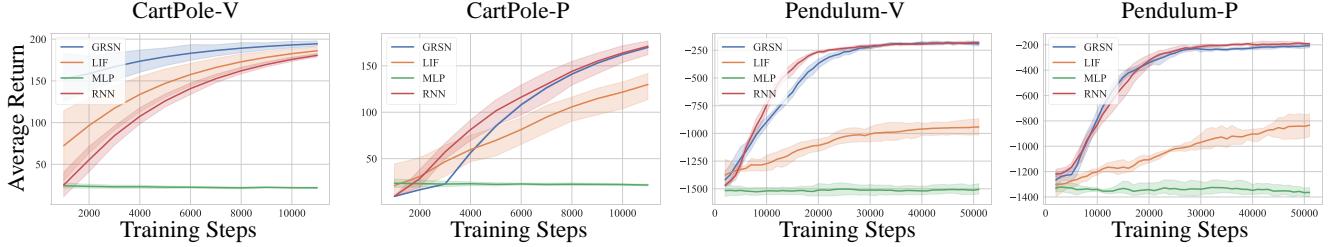
Figure 5: Results of classic partial observable control tasks. The shaded region represents a standard deviation of average evaluation over five different seeds, and the learning curves are smoothed for visual clarity.

| Enviroment | Pendulum-P | | | | SMAC-8m | | | |
|---|---|---|---|---|---|---|---|---|
| **Paradigm** | **w/ TAP** | | **w/o TAP** [†] | | **w/ TAP** | | **w/o TAP** [†] | |
| | **LIF** | **GRSN** | **LIF** | **GRSN** | **LIF** | **GRSN** | **LIF** | **GRSN** |
| **Mean Reward ↑** | -824.8 | **-195.8** | -1069.9 | -235.4 | 15.3 | **19.9** | 16.9 | 19.7 |
| **Standard Deviation ↓** | 178.80 | **43.99** | 65.85 | 44.16 | 1.59 | **0.02** | 1.38 | 0.31 |
| **Test Win Rate ↑ (%)** | - | - | - | - | 50.4 | **99.4** | 64.6 | 96.5 |
| **Total Time Steps[*]↓** | **64** | **64** | 256 | 256 | **26** | **26** | 104 | 104 |
| **Training Time ↓ (h)** | 3.55 | **3.42** | 7.26 | 10.18 | 12.91 | **12.76** | 19.52 | 18.57 |

[†] For all experiments without TAP, SNNs use repeated input and rate coding, and the time step is set to $T = 4$.
[*] In the SMAC environment, the average episode length is used instead of time steps.

Table 3: Ablation Study

## 5.4 Energy Consumption Estimation

SNNs use discrete spikes for information transmission, so their energy consumption is estimated differently than that of ANNs. We follow the convention of the neuromorphic computing community [Wang *et al.*, 2023; Qin *et al.*, 2023b; Wu *et al.*, 2023] by counting the total synaptic operations (SOP) to estimate the energy consumption of SNN models. Specifically, the SOP for SNNs correlates with the neurons' firing rate, fan-out $f_{out}$, and time window size $T$:

$$SOP = \sum_{t=1}^{T} \sum_{l=1}^{L-1} \sum_{j=1}^{N^l} f_{out}^l[j] \boldsymbol{o}_t^l[j] \tag{19}$$

where $L$ is the total number of layers and $N^l$ is the number of neurons in layer $l$, $f_{out}$ is the number of outgoing connections, and $\boldsymbol{o}_t^l[j]$ denotes the output spike of the $j$-th neuron in layer $l$ at the $t$-th time step.

For ANNs, SOP is only related to network structure. It is defined as:

$$SOP = \sum_{l=1}^{L} f_{in}^l N^l \tag{20}$$

where $f_{in}^l$ represents the number of incoming connections to each neuron in layer $l$, $L$ and $N^l$ are consistent with their definitions in SNNs.

SNNs use efficient addition operations, while ANNs use more expensive multiply-accumulate operations. According to [Han *et al.*, 2015], we measure 32-bit floating-point addition operations by $0.9pJ$ per operation and 32-bit floating-point multiply-accumulate operations by $4.6pJ$ per operation.

| Net. Env. | GRU (K$p$J) | | GRSN (K$p$J) | |
|---|---|---|---|---|
| | **MLP** | **RNN** | **MLP** | **SNN** |
| PO-Ctrl. | 536.06 | 628.82 | **498.34** | 7.92 |
| SMAC | **50.97** | 114.82 | 96.67 | **0.03** |
| Saved | - | | $\sim$ **49.1%** | |

Table 4: Energy Consumption Estimation

Here, we randomly select 1024 samples to estimate the average SOP and energy consumption for SNNs. Table 4 shows the results of different networks in different environments, and it can be seen that SNN can reduce energy consumption by up to about 3800 times compared to RNNs. The total energy consumption is also reduced by about 50%, which effectively reduces the agent's resource requirements.

## 6 Conclusion

In this paper, we introduce the temporal mismatch problem in SRL algorithms and propose the temporal alignment paradigm (TAP) and gated recurrent spiking neurons (GRSN) to solve this problem. Specifically, TAP allows the single-step update of spiking neurons to correspond to the single-step decision-making of MDP to solve the mismatch problem and significantly reduces time steps. On the other hand, GRSN enhances the temporal correlation of spiking neurons, thereby improving their ability to capture historical information and compensating for the performance degradation caused by time step reduction. Extensive experiments show that TAP can reduce the time steps by several times and half the

training time, which effectively solves the temporal mismatch problem. At the same time, GRSN can also improve the performance of spiking neurons, making SNNs equal to or even exceed the traditional RNNs. Energy estimation proves that our method can reduce the energy consumption by about 50%, which provides a new algorithm option for the control task in the energy-constrained scenario.

# References

[Barto *et al.*, 1983] Andrew G. Barto, Richard S. Sutton, and Charles W. Anderson. Neuronlike adaptive elements that can solve difficult learning control problems. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-13(5):834–846, 1983.

[Bellman, 1966] Richard Bellman. Dynamic programming. *Science*, 153(3731):34–37, 1966.

[Bohté *et al.*, 2002] Sander M. Bohté, Joost N. Kok, and Johannes A. La Poutré. Error-backpropagation in temporally encoded networks of spiking neurons. *Neurocomputing*, 48:17–37, 2002.

[Brockman *et al.*, 2016] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym, 2016.

[Chen *et al.*, 2022] Ding Chen, Peixi Peng, Tiejun Huang, and Yonghong Tian. Deep reinforcement learning with spiking q-learning. *CoRR*, abs/2201.09754, 2022.

[Chung *et al.*, 2014] Junyoung Chung, Çaglar Gülçehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *CoRR*, abs/1412.3555, 2014.

[Fang *et al.*, 2023] Wei Fang, Yanqi Chen, Jianhao Ding, Zhaofei Yu, Timothée Masquelier, Ding Chen, Liwei Huang, Huihui Zhou, Guoqi Li, and Yonghong Tian. Spikingjelly: An open-source machine learning infrastructure platform for spike-based intelligence. *Science Advances*, 9(40):eadi1480, 2023.

[Frémaux *et al.*, 2010] Nicolas Frémaux, Henning Sprekeler, and Wulfram Gerstner. Functional requirements for reward-modulated spike-timing-dependent plasticity. *Journal of Neuroscience*, 30:13326–13337, 2010.

[Fujimoto *et al.*, 2018] Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *Proceedings of the 35th International Conference on Machine Learning*, pages 1582–1591, 2018.

[Haarnoja *et al.*, 2018] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the 35th International Conference on Machine Learning*, pages 1856–1865, 2018.

[Han *et al.*, 2015] Song Han, Jeff Pool, John Tran, and William J. Dally. Learning both weights and connections for efficient neural networks. *CoRR*, abs/1506.02626, 2015.

[Han *et al.*, 2020] Dongqi Han, Kenji Doya, and Jun Tani. Variational recurrent models for solving partially observable control tasks. In *Proceedings of the 8th International Conference on Learning Representations*, 2020.

[Hu *et al.*, 2021] Jian Hu, Siyang Jiang, Seth Austin Harding, Haibin Wu, and Shih wei Liao. Rethinking the implementation tricks and monotonicity constraint in cooperative multi-agent reinforcement learning. *CoRR*, abs/2102.03479, 2021.

[Kingma and Ba, 2015] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of the 3rd International Conference on Learning Representations*, 2015.

[Ledinauskas *et al.*, 2020] Eimantas Ledinauskas, Julius Ruseckas, Alfonsas Jursenas, and Giedrius Burachas. Training deep spiking neural networks. *CoRR*, abs/2006.04436, 2020.

[Lian *et al.*, 2023] Shuang Lian, Jiangrong Shen, Qianhui Liu, Ziming Wang, Rui Yan, and Huajin Tang. Learnable surrogate gradient for direct training spiking neural networks. In *Proceedings of the 32nd International Joint Conference on Artificial Intelligence*, pages 3002–3010, 2023.

[Liu *et al.*, 2022] 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*, pages 1–12, 2022.

[Liu *et al.*, 2023] Qian Liu, Lifan Long, Hong Peng, Jun Wang, Qian Yang, Xiaoxiao Song, Agustín Riscos-Núñez, and Mario J. Pérez-Jiménez. Gated spiking neural P systems for time series forecasting. *IEEE Transactions on Neural Networks and Learning Systems*, 34:6227–6236, 2023.

[Lotfi-Rezaabad and Vishwanath, 2020] Ali Lotfi-Rezaabad and Sriram Vishwanath. Long short-term memory spiking networks and their applications. In *Proceedings of the International Conference on Neuromorphic Systems*, pages 3:1–3:9, 2020.

[Neftci *et al.*, 2019] Emre O. Neftci, Hesham Mostafa, and Friedemann Zenke. Surrogate gradient learning in spiking neural networks: Bringing the power of gradient-based optimization to spiking neural networks. *IEEE Signal Processing Magazine*, 36:51–63, 2019.

[Ni *et al.*, 2022] Tianwei Ni, Benjamin Eysenbach, and Ruslan Salakhutdinov. Recurrent model-free RL can be a strong baseline for many pomdps. In *Proceedings of the 39th International Conference on Machine Learning*, pages 16691–16723, 2022.

[Patel *et al.*, 2019] 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.

[Pei *et al.*, 2019] Jing Pei, Lei Deng, Sen Song, Mingguo Zhao, Youhui Zhang, Shuang Wu, Guanrui Wang, Zhe Zou, Zhenzhi Wu, Wei He, Feng Chen, Ning Deng, Si Wu, Yu Wang, Yujie Wu, Zheyu Yang, Cheng Ma, Guoqi Li, Wentao Han, Huanglong Li, Huaqiang Wu, Rong Zhao, Yuan Xie, and Luping Shi. Towards artificial general intelligence with hybrid tianjic chip architecture. *Nature*, 572:106–111, 2019.

[Ponghiran and Roy, 2022] Wachirawit Ponghiran and Kaushik Roy. Spiking neural networks with improved inherent recurrence dynamics for sequential learning. In *Proceedings of the 36th AAAI Conference on Artificial Intelligence*, pages 8001–8008, 2022.

[Qin *et al.*, 2023a] Lang Qin, Ziming Wang, Rui Yan, and Huajin Tang. Attention-based deep spiking neural networks for temporal credit assignment problems. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–11, 2023.

[Qin *et al.*, 2023b] Lang Qin, Rui Yan, and Huajin Tang. A low latency adaptive coding spike framework for deep reinforcement learning. In *Proceedings of the 32nd International Joint Conference on Artificial Intelligence*, pages 3049–3057, 2023.

[Rao *et al.*, 2022] Arjun Rao, Philipp Plank, Andreas Wild, and Wolfgang Maass. A long short-term memory for AI applications in spike-based neuromorphic hardware. *Nature Machine Intelligence*, 4:467–479, 2022.

[Rashid *et al.*, 2018] Tabish Rashid, Mikayel Samvelyan, Christian Schröder de Witt, Gregory Farquhar, Jakob N. Foerster, and Shimon Whiteson. QMIX: monotonic value function factorisation for deep multi-agent reinforcement learning. In *Proceedings of the 35th International Conference on Machine Learning*, pages 4292–4301, 2018.

[Rathi and Roy, 2023] Nitin Rathi and Kaushik Roy. DIET-SNN: A low-latency spiking neural network with direct input encoding and leakage and threshold optimization. *IEEE Transactions on Neural Networks and Learning Systems*, 34(6):3174–3182, 2023.

[Roy *et al.*, 2019] Kaushik Roy, Akhilesh Jaiswal, and Priyadarshini Panda. Towards spike-based machine intelligence with neuromorphic computing. *Nature*, 575:607–617, 2019.

[Samvelyan *et al.*, 2019] Mikayel Samvelyan, Tabish Rashid, Christian Schröder de Witt, Gregory Farquhar, Nantas Nardelli, Tim G. J. Rudner, Chia-Man Hung, Philip H. S. Torr, Jakob N. Foerster, and Shimon Whiteson. The starcraft multi-agent challenge. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*, pages 2186–2188, 2019.

[Seung, 2003] H Sebastian Seung. Learning in spiking neural networks by reinforcement of stochastic synaptic transmission. *Neuron*, 40:1063–1073, 2003.

[Sun *et al.*, 2022] 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, 2022.

[Sutton and Barto, 2018] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

[Tan *et al.*, 2021] 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 35th AAAI Conference on Artificial Intelligence*, pages 9816–9824, 2021.

[Tang *et al.*, 2020a] 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*, pages 6090–6097, 2020.

[Tang *et al.*, 2020b] Guangzhi Tang, Neelesh Kumar, Raymond Yoo, and Konstantinos P. Michmizos. Deep reinforcement learning with population-coded spiking neural network for continuous control. In *Proceedings of the 2020 Conference on Robot Learning*, pages 2016–2029, 2020.

[Urbanczik and Senn, 2009] Robert Urbanczik and Walter Senn. Reinforcement learning in populations of spiking neurons. *Nature Neuroscience*, 3:250–252, 2009.

[Wang *et al.*, 2023] Ziming Wang, Runhao Jiang, Shuang Lian, Rui Yan, and Huajin Tang. Adaptive smoothing gradient learning for spiking neural networks. In *Proceedings of the 40th International Conference on Machine Learning*, pages 35798–35816, 2023.

[Williams, 1992] Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8:229–256, 1992.

[Wu *et al.*, 2018] 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, 2018.

[Wu *et al.*, 2023] Jibin Wu, Yansong Chua, Malu Zhang, Guoqi Li, Haizhou Li, and Kay Chen Tan. A tandem learning rule for effective training and rapid inference of deep spiking neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 34(1):446–460, 2023.

[Yin *et al.*, 2020] Bojian Yin, Federico Corradi, and Sander M. Bohté. Effective and efficient computation with multiple-timescale spiking recurrent neural networks. In *Proceedings of the International Conference on Neuromorphic Systems*, pages 1:1–1:8, 2020.

[Zhang *et al.*, 2021] Ling Zhang, Jian Cao, Yuan Zhang, Bohan Zhou, and Shuo Feng. Distilling neuron spike with high temperature in reinforcement learning agents. *CoRR*, abs/2108.10078, 2021.

[Zhang *et al.*, 2022] Duzhen Zhang, Tielin Zhang, Shuncheng Jia, and Bo Xu. Multiscale dynamic coding improved spiking actor network for reinforcement learning. In *Proceedings of the 36th AAAI Conference on Artificial Intelligence*, pages 59–67, 2022.

## A   Background

In reinforcement learning (RL), the process of interacting with environments and changing the state through different actions is often described by Markov Decision Process (MDP) $(\mathcal{S}, \mathcal{A}, p_M, r, \gamma)$, with state space $\mathcal{S}$, action space $\mathcal{A}$, discount factor $\gamma$, and transition dynamics $p_M(s'|s, a)$. At every step, the agent change its state from $s$ to $s'$ by perform action $a$ and receives a reward $r(s, a, s')$. The goal of RL is to allow the agent to maximize the expectation of accumulating rewards $R_t = \sum_{i=t+1}^{\infty} \gamma^i r(s_i, a_i, s_{i+1})$.

The agent selects actions with respect to a policy $\pi : \mathcal{S} \to \mathcal{A}$. The corresponding value function $Q^\pi(s, a) = \mathbb{E}_\pi[R_t|s, a]$, the expected accumulate reward of agent for taking action $a$ in state $s$. Rewriting the value function $Q^\pi(s, a)$ into an iterative form, we can get the Bellman equation:

$$Q^\pi(s, a) = \mathbb{E}_\pi[r + \gamma Q^\pi(s', a')] \tag{21}$$

Then we can define the optimal value function:

$$Q^*(s, a) = \max_\pi Q^\pi(s, a) \tag{22}$$

The policy $\pi^*$ that achieves the optimal value function is the optimal policy. It can be obtained through greedy action choices. For high-dimensional control tasks, the value function can be fitted by deep neural networks (DNNs), which is the core idea of DRL.

## B   Proof

According to the proof of Theorem 1, the third term in Eq 9 ($\frac{\partial \boldsymbol{u}_t^l}{\partial \boldsymbol{W}^l}$) determines the direction of gradient flow. $\frac{\partial \boldsymbol{u}_t^l}{\partial \boldsymbol{W}^l}$ can be calculated as:

$$\frac{\partial \boldsymbol{u}_t^l}{\partial \boldsymbol{W}^l} = (1 - \beta) \sum_{i=1}^{t} \beta^i \prod_{j=1}^{i} \delta_{t-j} \frac{\partial \boldsymbol{c}_{t-i}^l}{\partial \boldsymbol{W}^l} + \prod_{k=1}^{t} \beta \delta_{t-k} + (1 - \beta) \frac{\partial \boldsymbol{c}_t^l}{\partial \boldsymbol{W}^l} \tag{23}$$

In Eq 23, both the first and third terms contain $\frac{\partial \boldsymbol{c}}{\partial \boldsymbol{W}}$, which is consistent with the assumption that the gradient mainly flows along the current. Therefore, we only need to analyze the second term in Eq 23.

### B.1   Hard Reset

When the neuron model adopted hard reset, the $\delta_t$ can be calculated as:

$$\delta_t = \frac{\partial \hat{\boldsymbol{u}}_t^l}{\partial \boldsymbol{u}_t^l} = \frac{\partial [u_r \boldsymbol{o}_t^l + (1 - \boldsymbol{o}_t^l) \boldsymbol{u}_t^l]}{\partial \boldsymbol{u}_t^l} = [(u_r - \boldsymbol{u}_t^l) h'(\boldsymbol{u}_t^l - \vartheta) + 1 - \boldsymbol{o}_t^l] \tag{24}$$

so the second term in Eq 23 can be simplified as:

$$\prod_{k=1}^{t} \beta \delta_{t-k} = \prod_{k=1}^{t} \beta [(u_r - \boldsymbol{u}_{t-k}^l) h'(\boldsymbol{u}_{t-k}^l - \vartheta) + 1 - \boldsymbol{o}_{t-k}^l] \tag{25}$$

where $h'(x)$ is the surrogate gradient function:

$$h'(x) = \frac{\alpha}{2[1 + (\frac{\pi}{2} \alpha x)^2]} \tag{26}$$

Based on the experimental parameter settings $\boldsymbol{u}_r = 0, \vartheta = 1, \alpha = 2, \beta = 0.5$, we can simplify the multiplication term:

$$\begin{aligned}
&\beta [(u_r - \boldsymbol{u}_t^l) h'(\boldsymbol{u}_t^l - \vartheta) + 1 - \boldsymbol{o}_t^l] \\
&= \frac{1}{2} [1 - \boldsymbol{o}_t^l - \boldsymbol{u}_t^l h'(\boldsymbol{u}_t^l - 1)] \\
&= \frac{1}{2} [1 - \boldsymbol{o}_t^l - \frac{\boldsymbol{u}_t^l}{1 + \pi^2 (\boldsymbol{u}_t^l - 1)^2}]
\end{aligned} \tag{27}$$

When $\boldsymbol{u}_t^l \geq \vartheta = 1$, neurons will emit spikes ($\boldsymbol{o}_t^l = 1$). Eq 27 can be further simplified as:

$$\begin{aligned}
&\frac{1}{2} [1 - \boldsymbol{o}_t^l - \frac{\boldsymbol{u}_t^l}{1 + \pi^2 (\boldsymbol{u}_t^l - 1)^2}] \\
&= -\frac{1}{2} \frac{\boldsymbol{u}_t^l}{1 + \pi^2 (\boldsymbol{u}_t^l - 1)^2}
\end{aligned} \tag{28}$$

Let $f(x) = -\frac{1}{2}\frac{x}{1+\pi^2(x-1)^2}$, then we can analyze the monotonicity of this item.

$$f(x) = -\frac{1}{2}\frac{x}{1+\pi^2(x-1)^2}$$
$$f'(x) = -\frac{1}{2}\frac{1+\pi^2(x-1)^2 - 2\pi^2(x-1)x}{[1+\pi^2(x-1)^2]^2} \tag{29}$$
$$= -\frac{1}{2}\frac{1+\pi^2-\pi^2 x^2}{[1+\pi^2(x-1)^2]^2}$$

Let $f'(x) = 0$, then:

$$x = \pm\sqrt{1+\frac{1}{\pi^2}}$$

$$\Rightarrow max[f(x)] = f(-\sqrt{1+\frac{1}{\pi^2}}) = \frac{\frac{1}{2}\sqrt{1+\frac{1}{\pi^2}}}{2+2\pi^2+2\pi^2\sqrt{1+\frac{1}{\pi^2}}} \approx 0.0124 \tag{30}$$

$$\Rightarrow min[f(x)] = f(\sqrt{1+\frac{1}{\pi^2}}) = -\frac{\frac{1}{2}\sqrt{1+\frac{1}{\pi^2}}}{2+2\pi^2-2\pi^2\sqrt{1+\frac{1}{\pi^2}}} \approx -0.5124$$

According to Eq 30, we can find $|f(x)| < 1$, which means $|\beta\delta_t| < 1$. So the term $\prod_{k=1}^{t}\beta\delta_{t-k} \to 0$.
In another case, when $\boldsymbol{u}_t^l < \vartheta = 1$, neurons will remain silent ($\boldsymbol{o}_t^l = 0$). Then Eq 27 can be simplified as:

$$\frac{1}{2}[1 - \boldsymbol{o}_t^l - \frac{\boldsymbol{u}_t^l}{1+\pi^2(\boldsymbol{u}_t^l - 1)^2}]$$
$$= \frac{1}{2}[1 - \frac{\boldsymbol{u}_t^l}{1+\pi^2(\boldsymbol{u}_t^l - 1)^2}] \tag{31}$$

Let $g(x) = \frac{1}{2}(1 - \frac{x}{1+\pi^2(x-1)^2})$, then we can find:

$$g(x) = f(x) + \frac{1}{2} \tag{32}$$

Similarly, we can calculate the maximum and minimum value of $g(x)$:

$$max[g(x)] = max[f(x)] + \frac{1}{2} \approx 0.5124$$
$$\tag{33}$$
$$min[g(x)] = min[f(x)] + \frac{1}{2} \approx -0.0124$$

This also satisfied $|g(x)| < 1$, which means $\prod_{k=1}^{t}\beta\delta_{t-k} \to 0$.
Based on Eq 25-33, we can conclude that the second term $\prod_{k=1}^{t}\beta\delta_{t-k}$ in Eq 23 tends to 0 when total time steps $T$ is large. So the Eq 23 can be simplified as:

$$\frac{\partial\boldsymbol{u}_t^l}{\partial\boldsymbol{W}^l} = (1-\beta)\sum_{i=1}^{t}\beta^i\prod_{j=1}^{i}\delta_{t-j}\frac{\partial\boldsymbol{c}_{t-i}^l}{\partial\boldsymbol{W}^l} + (1-\beta)\frac{\partial\boldsymbol{c}_t^l}{\partial\boldsymbol{W}^l} \tag{34}$$

which means the gradient of loss function $\frac{\partial\mathcal{L}}{\partial\boldsymbol{W}^l}$ mainly depends on current $\frac{\partial\boldsymbol{c}}{\partial\boldsymbol{W}^l}$.

## B.2  Soft Reset

When the neuron model adopted soft reset, the $\delta_t$ should be calculated as:

$$\delta_t = \frac{\partial\hat{\boldsymbol{u}}_t^l}{\partial\boldsymbol{u}_t^l} = \frac{\partial(\boldsymbol{u}_t^l - \vartheta\boldsymbol{o}_t^l)}{\partial\boldsymbol{u}_t^l} = 1 - \vartheta h'(\boldsymbol{u}_t^l - \vartheta) \tag{35}$$

then the second term in Eq 23 can be simplified as:

$$\prod_{k=1}^{t}\beta\delta_{t-k} = \prod_{k=1}^{t}\beta[1 - \vartheta h'(\boldsymbol{u}_{t-k}^l - \vartheta)] \tag{36}$$
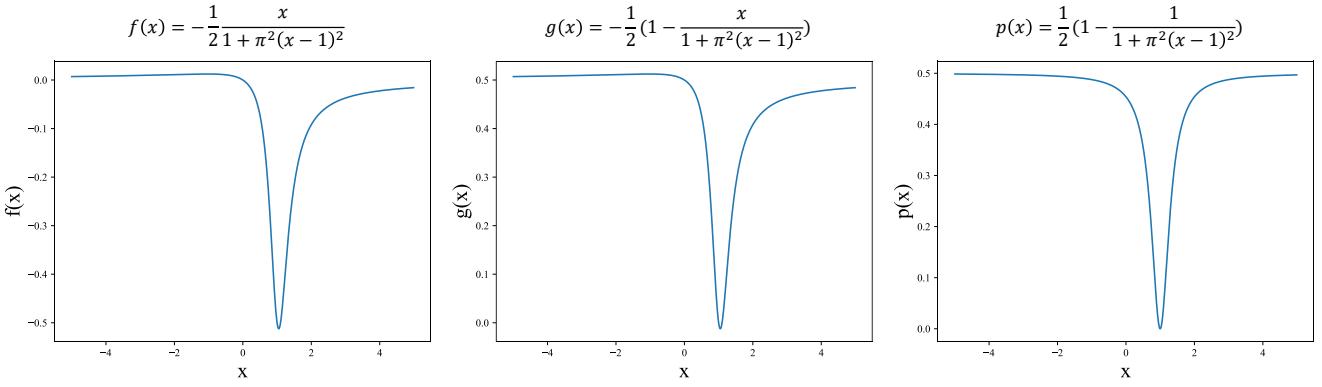
Figure 6: Graph of function $f(x)$, $g(x)$ and $p(x)$.

where $h'(x)$ is the same as Eq 26.

Similar to the hard reset mode, we bring the values $(\boldsymbol{u}_r = 0, \vartheta = 1, \alpha = 2, \beta = 0.5)$ of each parameter into the equation:

$$\beta[1 - \vartheta h'(\boldsymbol{u}_t^l - \vartheta)] = \frac{1}{2}[1 - \frac{1}{1 + \pi^2(\boldsymbol{u}_t^l - 1)^2}] \tag{37}$$

Let $p(x) = \frac{1}{2}[1 - \frac{1}{1+\pi^2(x-1)^2}]$, then we can analyze this item:

$$\begin{aligned}
&\because \pi^2(x-1)^2 \geq 0 \\
&\Rightarrow 1 + \pi^2(x-1)^2 \geq 1 \\
&\Rightarrow 0 < \frac{1}{1 + \pi^2(x-1)^2} \leq 1 \\
&\Rightarrow 0 \leq 1 - \frac{1}{1 + \pi^2(x-1)^2} < 1 \\
&\Rightarrow 0 \leq p(x) < \frac{1}{2}
\end{aligned} \tag{38}$$

Obviously, this also satisfied $|p(x)| < 1$, so $\prod_{k=1}^{t} \beta\delta_{t-k} \to 0$.

### B.3 Curves of funtions

To further confirm the correctness of the monotonicity of functions, we visualize those three functions in Figure 6.

## C Experiment Details

### C.1 Experiment Settings

The experiments are all carried out on a workstation equipped with AMD Ryzen Threadripper 3970X CPU running at 2.90 GHz, 128 GB of RAM, and four NVIDIA Geforce RTX 3090 GPUs. We use Python 3.8.13 to process datasets and PyTorch 1.12.1 to implement deep SNNs. We also use the spikingjelly 0.0.0.0.14 [Fang *et al.*, 2023] framework to implement deep SNN networks and custom neurons. Table 5 shows various hyperparameters related to spiking neurons and spiking neural networks.

| Parameter | Meanings | Value |
|:---:|:---:|:---:|
| $\beta$ | Decay rate of membrane voltage | 0.5 |
| $\alpha$ | Width of surrogate function | 2 |
| $\vartheta$ | Potential threshold | 1 |
| $T$ | Default length of time window | 4 |
| $\boldsymbol{u}_r$ | Reset potential | 0 |

Table 5: Value of SNN parameters

## C.2 POMDP

We used the initial versions (-v0) of *CartPole* and *Pendulum* environments for testing and packaged the environment with Openai gym [Brockman *et al.*, 2016]. Table 6 and Table 7 shows all the hyperparameters in our POMDP experiments, and Figure 7 shows the overview of network structure. The setting of each parameter and the design of the network structure all come from [Ni *et al.*, 2022].

| Parameter | *Pendulum* | *CartPole* |
|---|---|---|
| Training steps | 50K | 10K |
| Meta-training steps | 250 | 50 |
| Max environment steps | 200 | 200 |
| Evaluate interval steps | 5 | 1 |
| Agent inputs | Observation/Action/Reward | Observation/Action/Reward |

Table 6: Hyperparameters of partial observable environments

| Meanings | Parameter | Value |
|---|---|---|
| Hidden layer size | Observation embedder | [32] |
| | Action embedder | [8] |
| | Reward embedder | [8] |
| | MLP | [128,128] |
| RL parameters | Optimizer | Adam [Kingma and Ba, 2015] |
| | Learning rate | 3e-4 |
| | Discount factor $\gamma$ | 0.9 |
| | Smoothing coef $\tau$ | 0.005 |
| | Replay buffer size | 1e6 |
| | Batch size | 32 |
| RNN | Sequence length | 64 |
| | RNN hidden size | [128] |

Table 7: Hyperparameters of algorithms and networks

For each set of experiments, we will save the results of the model testing during the training process, and ultimately choose the mean of the last five test results during the training process as the final result of this experiment. Then, we will select five different random seeds to independently repeat the experiment for five times, and the mean and standard deviation of the final five results will be the final result of this set of experiments.

## C.3 MARL

Due to the weak parameter settings of the original QMIX algorithm, we referred to the paper [Hu *et al.*, 2021] for parameter adjustments, and the specific parameters are shown in the Table 8. For each independent experiment, we use the average of the last fifty testing results to represent its performance. The final results of the algorithm in different environments are the average performance of five independent experiments.

| Parameter | Value |
|---|---|
| Optimizer | Adam |
| Batch size | 128 |
| Q($\lambda$) | 0.6 |
| Attention heads | - |
| Mixing network size | 41K |
| $\epsilon$ Anneal steps | 100K |
| Rollout processes | 8 |

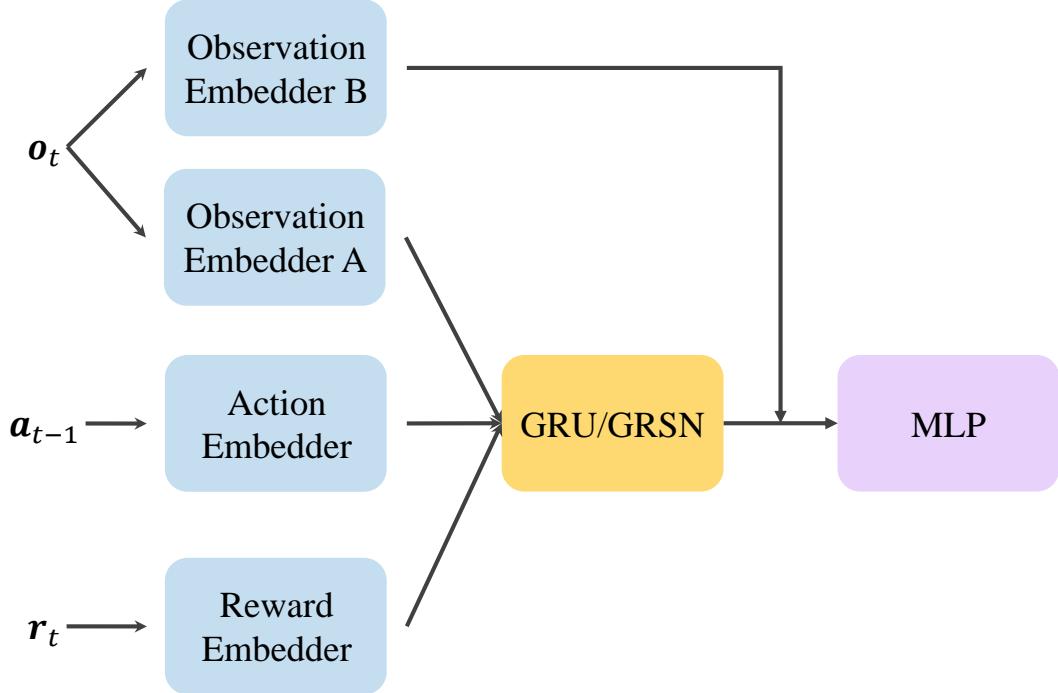Table 8: Hyperparameters of our QMIX algorithm

Figure 7: The network architecture we used in POMDP experiments.

## C.4   Energy Estimation

The main text only displays the final estimation results, and here we present the complete data in Table 9.

| Enviroment | | GRU | | | GRSN | | | Mixer |
|---|---|---|---|---|---|---|---|---|
| | | MLP (K) | RNN (K) | Energy (K$pJ$) | MLP (K) | SNN (K) | Energy (K$pJ$) | |
| PO-Ctrl. | *Pendulum-P* | 116.39 | | 1164.22 | 108.19 | 9.73 | **506.44** | |
| | *Pendulum-V* | 116.24 | 136.70 | 1163.56 | 108.05 | 11.16 | **507.07** | - |
| | *CartPole-P* | 116.74 | | 1165.85 | 108.55 | 6.98 | **505.91** | |
| | *CartPole-V* | 116.74 | | 1165.85 | 108.55 | 7.33 | **505.60** | |
| MARL | *8m* | 7.57 | | 149.62 | 14.98 | 0.01 | **68.90** | 51.20 |
| | *2s3z* | 6.99 | | 146.97 | 22.85 | 0.01 | **105.11** | 35.75 |
| | *8m_vs_9m* | 8.02 | 24.96 | 151.69 | 15.36 | 0.06 | **70.71** | 53.31 |
| | *3s_vs_5z* | 4.81 | | 136.95 | 12.42 | 0.03 | **57.15** | 21.60 |
| | *27m_vs_30m* | 24.74 | | 228.62 | 30.72 | 0.07 | **141.38** | 283.11 |
| | *MMM2* | 14.35 | | 180.84 | 29.76 | 0.01 | **136.90** | 84.93 |

Table 9: Energy Estimation