

---

# NBSP: A Neuron-Level Framework for Balancing Stability and Plasticity in Deep Reinforcement Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 In contrast to the human ability to continuously acquire knowledge, agents struggle  
2 with the stability-plasticity dilemma in deep reinforcement learning (DRL), which  
3 refers to the trade-off between retaining existing skills (stability) and learning new  
4 knowledge (plasticity). Current methods focus on balancing these two aspects  
5 at the network level, lacking sufficient differentiation and fine-grained control  
6 of individual neurons. To overcome this limitation, we propose Neuron-level  
7 Balance between Stability and Plasticity (NBSP) method, by taking inspiration  
8 from the observation that specific neurons are strongly relevant to task-relevant  
9 skills. Specifically, NBSP first (1) defines and identifies RL skill neurons that  
10 are crucial for knowledge retention through a goal-oriented method, and then (2)  
11 introduces a framework by employing adaptive gradient masking and experience  
12 replay techniques targeting these neurons to preserve the encoded existing skills  
13 while enabling adaptation to new tasks. Numerous experimental results on the Meta-  
14 World and Atari benchmarks demonstrate that NBSP significantly outperforms  
15 existing approaches in balancing stability and plasticity.

## 16 1 Introduction

17 **Deep reinforcement learning (DRL)** has shown exceptional capabilities across a range of complex  
18 scenarios, such as gaming (Mnih et al., 2013), robotic manipulation (Andrychowicz et al., 2020), and  
19 autonomous driving (Kiran et al., 2021). However, most RL research focuses on agents that learn to  
20 solve individual problems rather than learn a sequence of tasks continually. Ideally, the agent must  
21 maintain its performance on previously learned tasks, referred to as **stability** (McCloskey & Cohen,  
22 1989), while simultaneously adapting to new tasks, known as **plasticity** (Carpenter & Grossberg,  
23 1987). However, it has been revealed that emphasizing stability may hinder the ability of agents to  
24 learn new knowledge (Nikishin et al., 2022a; Abbas et al., 2023), whereas excessive plasticity can  
25 lead to catastrophic forgetting of previously acquired knowledge (Goodfellow et al., 2015; Atkinson  
26 et al., 2021b), a challenge known as the **stability-plasticity dilemma** (eMermillod et al., 2013),  
27 which remains a fundamental and under-explored problem and is the main focus of our work.

28 Existing methods to strike a balance between stability and plasticity generally fall into three categories,  
29 i.e. (1) **regularization-based methods** (Kirkpatrick et al., 2017; Kumar et al., 2023), which apply  
30 penalties to parameter changes to mitigate forgetting while acquiring new knowledge; (2) **replay-**  
31 **based methods** (Ahn et al., 2024), which leverage past experiences to consolidate knowledge; and  
32 (3) **modularity-based methods** (Kim et al., 2023; Anand & Precup, 2024), which seek to decouple  
33 stability and plasticity or isolate different components for different tasks. Despite their contributions,  
34 these methods suffer from three key limitations: (1) They primarily operate at the network level, yet  
35 their ultimate impact manifests at the level of individual neurons. However, these methods fail to

36 differentiate and fine-grained control neurons based on their specific roles. Therefore, identifying  
 37 and effectively utilizing task-relevant neurons remains both critical and under-explored. (2) These  
 38 studies are primarily conducted within the paradigm of continual learning, thus overlooking the  
 39 unique characteristics intrinsic to DRL. (3) These approaches could sometimes unnecessarily inflate  
 40 model parameters, thereby introducing unwarranted complexity (Bai et al., 2023).

41 By analyzing the activations of neurons in the DRL network, we observe that the activations of certain neurons  
 42 are strongly correlated with the task goal. For instance, Figure 1 illustrates the activation distribution of a specific  
 43 neuron in the network following training on the drawer-open task from the Meta-World benchmark Yu et al. (2020).  
 44 Activation of the neuron serves as a reliable predictor of whether the task is successful. Higher activation levels  
 45 correspond to an increased likelihood of completing the task successfully, indicating that this neuron encodes a  
 46 critical skill essential for the task. Consequently, it plays a pivotal role in retaining task-specific memory.  
 47  
 48  
 49  
 50  
 51  
 52

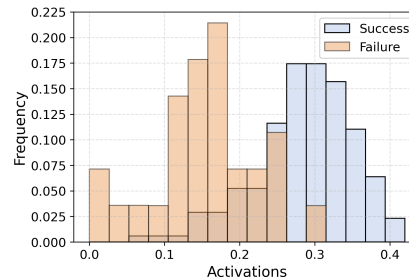


Figure 1: Distribution histogram of the activation of a neuron, categorized based on whether the drawer-open task was successfully completed or not.

53 Motivated by the aforementioned observations, we propose **Neuron-level Balance between Stability and Plasticity (NBSP)**, a novel DRL framework that operates at  
 54 the level of neurons to tackle the stability-plasticity dilemma. In particular, (1) we first introduce **RL skill neurons**, which encode critical skills necessary for knowledge retention. While skill neurons  
 55 have been investigated and successfully exploited in various domains, such as pre-trained language models (Wang et al., 2022) and neural machine translation (Bau et al., 2018), skill neurons are still  
 56 much less explored in DRL. We bridge this research gap by proposing a goal-oriented strategy for identifying RL skill neurons. (2) We then apply **gradient masking** according to the scores of these  
 57 neurons, ensuring that the encoded knowledge of prior skills is preserved while allowing fine-tuning during subsequent training. Meanwhile, the other neurons retain the ability to learn new tasks. (3)  
 58 Additionally, we incorporate **experience replay** to periodically revisit the past experience to reinforce stability, preventing excessive drift from previous knowledge. Integrally, NBSP offers three key  
 59 advantages compared with previous methods: (1) The neuron-level processing enables finer control and greater flexibility, addressing the stability-plasticity trade-off at the most fundamental level of the  
 60 network. (2) The goal-oriented approach for identifying RL skill neurons is specifically tailored to DRL. (3) This framework is simple, avoiding complex network designs or additional parameters.  
 61  
 62  
 63  
 64  
 65  
 66  
 67  
 68  
 69

70 We conduct experiments on the **Meta-World** (Yu et al., 2020) and **Atari** (Mnih et al., 2013) benchmarks to evaluate the effectiveness of NBSP. Our results show that NBSP outperforms existing  
 71 methods in balancing stability and plasticity, enabling effective learning of new tasks while preserving knowledge from previous tasks. Additionally, we perform extensive ablation studies to investigate  
 72 the contribution of different components within NBSP. Specially, we analyze the DRL agents by dissecting the performance of the two critical modules, i.e., the actor and the critic. Our findings  
 73 reveal that (1) addressing both the actor and critic networks yields the best performance, and (2) the critic plays a more critical role in achieving this balance due to the differences in their inherent  
 74 training mechanisms. In summary, our key contributions include:  
 75  
 76  
 77  
 78

- 79 • We are the first to introduce the concept of RL skill neurons which encode skills of the task, essential for knowledge retention, and propose a goal-oriented strategy specifically tailored to  
 80 DRL for identification.
- 81
- 82 • We tackle the stability-plasticity dilemma in DRL from the perspective of RL skill neurons, by employing gradient masking and experience replay on these neurons, eliminating requirements of  
 83 complex network designs or additional parameters.
- 84
- 85 • We conduct extensive experiments on the Meta-World and Atari benchmarks to demonstrate the effectiveness of our method in balancing stability and plasticity.  
 86

## 87 2 Related Work

88 **Balance between stability and plasticity.** In DRL, addressing the stability-plasticity dilemma (Carpenter & Grossberg, 1988) has inspired various strategies. Stability-focused methods often utilize

90 replay techniques, such as A-GEM (Chaudhry et al., 2018b), using episodic memory to constrain  
 91 loss, and ClonEx-SAC (Wolczyk et al., 2022), enhancing performance through behavior cloning.  
 92 Pseudo-rehearsals from generative models further reduce storage requirements (Atkinson et al.,  
 93 2021a). Plasticity-focused methods aim to preserve network expressiveness, with solutions like CBP  
 94 (Dohare et al., 2024), resetting (Nikishin et al., 2022b), plasticity injection (Nikishin et al., 2024),  
 95 Reset & Distillation (Ahn et al., 2024), and CRelu (Abbas et al., 2023) to prevent activation collapse.  
 96 Modularity-based methods balance stability and plasticity by decoupling task-specific knowledge,  
 97 exemplified by soft modularity for routing networks (Yang et al., 2020), value function decomposition  
 98 (Anand & Precup, 2024), and compositional frameworks leveraging neural components (Mendez et al.,  
 99 2022). Methods such as CRelu and ClonEx-SAC focus on continual reinforcement learning (CRL),  
 100 but our study specifically targets the intrinsic balance between stability and plasticity, with other  
 101 factors such as task order controlled in a cycling task setting. Moreover, while most methods operate  
 102 at the network level, our approach explores neuron-level research, providing fine-grained control.

103 **Neuron-level research.** Recent research has shown that neuron sparsity often correlates with task-  
 104 specific performance (Xu et al., 2024), driving a growing focus on skill neurons to interpret network  
 105 behavior and tackle challenges across domains. For example, skill neurons have been used to  
 106 enhance transferability and efficiency in Transformers via pruning (Wang et al., 2022), and dormant  
 107 neurons have been recycled to improve training in DRL (Sokar et al., 2023). Other studies, such as  
 108 identifying Rosetta Neurons (Dravid et al., 2023) and language-specific neurons (Tang et al., 2024),  
 109 have advanced alignment and interpretability. However, neuron-level studies in DRL are still limited,  
 110 with methods like CoTASP (Yang et al., 2023) and PackNet (Mallya & Lazebnik, 2018) focusing  
 111 on task-specific sub-network selection via sparse prompts, pruning, and re-training. And NPC (Paik  
 112 et al., 2019) restricts important neurons to maintain stability. In contrast, our work identifies RL skill  
 113 neurons specific to DRL, balancing stability and plasticity with encoded task-relevant knowledge.

## 114 3 Methodology

115 In this section, we first introduce the terminology of RL skill neurons and then propose the Neuron-  
 116 level Balance between Stability and Plasticity (NBSP) method.

### 117 3.1 Problem Setup

118 We focus on the setting of sequential task learning without constraints on the time intervals between  
 119 tasks. In this setting, the agent is expected to perform all previously learned tasks after training,  
 120 without relying on task-specific signals. For instance, large models such as DeepSeek employ RL  
 121 to enhance their reasoning capabilities. However, different tasks, such as vision and mathematics,  
 122 demand distinct reasoning abilities. To first strengthen a specific type of reasoning and then generalize  
 123 to others, it is essential to strike a balance between stability and plasticity during sequential training.  
 124 Furthermore, in real-world applications, the enhanced model should be able to handle all tasks  
 125 without relying on explicit task signals. Let  $\tau \in \{\tau_1, \tau_2, \dots\}$  represent a sequence of task, each task  $\tau$   
 126 corresponds to a distinct Markov Decision Process (MDP)  $M^\tau = (S^\tau, A^\tau, P^\tau, R^\tau, \gamma^\tau)$ , where  $S^\tau$ ,  
 127  $A^\tau$ ,  $P^\tau$ ,  $R^\tau$  and  $\gamma^\tau$  denote the state space, action space, transition dynamics, reward function, and  
 128 discount factor, respectively. Instead of addressing a single MDP, the goal is to solve a sequence of  
 129 MDPs one by one using a universal policy  $\pi(a|s)$  and Q-function  $Q(s, a)$ . The primary challenge  
 130 lies in balancing plasticity, which refers to maximizing the discounted return of the current task, and  
 131 stability, which emphasizes the maximization of the expected discounted return averaged across all  
 132 previous tasks. This trade-off constitutes the core problem addressed in this work.

### 133 3.2 Identifying RL Skill Neurons

134 In this study, we make a key observation that the stability and plasticity of the agent network are  
 135 closely related to its expressive capabilities, which are significantly influenced by the behavior of  
 136 individual neurons. As evidenced in Molchanov et al. (2022), neuron expression determines how  
 137 information is propagated and processed, directly affecting the learning and knowledge retention  
 138 capabilities of the network. Therefore, understanding and controlling neuron behavior is at the most  
 139 fundamental level for striking a balance between stability and plasticity. On the one hand, when  
 140 neuron expression is stable and generalized, the agent network tends to exhibit high stability. On the  
 141 other hand, strong plasticity can be achieved given neuron expression is flexible and adaptable.

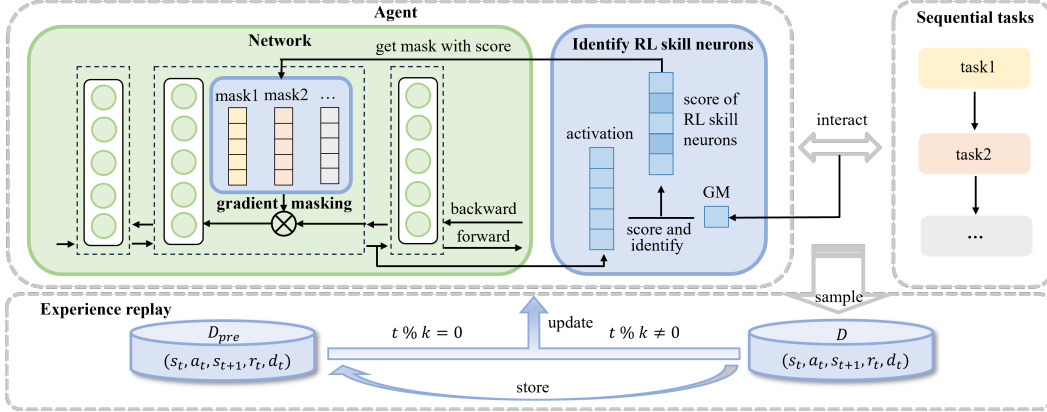


Figure 2: Framework of NBSP. The agent scores and identifies RL skill neurons for each task by measuring the activation in relation to the GM. While learning new tasks, the gradient of these neurons is masked adaptively based on their scores to preserve the encoded skills, while still allowing fine-tuning for new task learning. Additionally, a replay buffer is used to store a portion of the experiences from previous tasks, which is periodically sampled to update the agent.

142 Several works have demonstrated the multifaceted capabilities of neurons, such as the storage of  
 143 factual knowledge (Dai et al., 2022), the association with specific languages (Tang et al., 2024), and  
 144 the encoding of safety information (Chen et al., 2024). These specialized neurons, often referred as  
 145 skill neurons, have been shown to significantly contribute to network performance (Wang et al., 2022).  
 146 However, the potential of skill neurons in DRL remains largely under-explored. As illustrated in  
 147 Figure 1, activations of the specific neuron are strongly correlated with task success: higher activation  
 148 levels increase the likelihood of successful task completion, whereas lower levels are associated with  
 149 failure. *This indicates that the activations of these neurons significantly affect agent performance,*  
 150 *effectively encoding the critical skills required for the task. By preserving the activations of such*  
 151 *neurons, it becomes possible to retain the learned task-specific skills, thereby improving stability.*

152 In this work, we formally define these special neurons as **RL skill neurons**, which encode critical  
 153 skills, essential for knowledge retention in DRL. Furthermore, we propose a goal-oriented method  
 154 for the identification of these neurons. Unlike prior approaches that primarily focus on the inputs  
 155 triggering neuron activations (Bau et al., 2020; Gurnee & Tegmark, 2023), our method emphasizes  
 156 their impact on achieving ultimate goals, i.e. succeeding in finishing Meta-World tasks and attaining  
 157 high scores in Atari games, by comparing the activation patterns of the neurons that exhibit varying  
 158 performance levels. In Section 4.2, we empirically show the advantage of our goal-oriented method.

159 For a specific neuron  $\mathcal{N}$ , let  $a(\mathcal{N}, t)$  represent its activation at step  $t$ . In fully connected layers, each  
 160 output dimension corresponds to the activation of a specific neuron, whereas in convolution layers,  
 161 the average of each output channel represents the activation of a neuron. To quantify activation level  
 162 of a neuron  $\mathcal{N}$ , we define the **average activation** as:

$$\bar{a}(\mathcal{N}) = \frac{1}{T_{avg}} \sum_{t=1}^T a(\mathcal{N}, t), \quad (1)$$

163 where  $T_{avg}$  represents the average step. The activation level of the neuron can then be assessed by  
 164 comparing its current activation with the corresponding average activation.

165 To assess the performance of the agent at step  $t$ , we introduce the **Goal Metric (GM)**, denoted as  
 166  $q(t)$ . It serves as an evaluation metric for assessing the performance of the agent’s network, varying  
 167 based on the objective of the task. It is computed in an online manner during training. For instance,  
 168 on the Meta-World benchmark, the GM is typically binary, determined by whether the episode is  
 169 successful, which is computed at the end of each episode. In contrast, the GM is determined by the  
 170 cumulative return of the episode for the Atari benchmark. Additionally, we define the **average Goal**  
 171 **Metric (GM)** of the agent as follows, which serves as a baseline for evaluating the performance by  
 172 comparing it with the current GM.

$$\bar{q} = \frac{1}{T_{avg}} \sum_{t=1}^T q(t). \quad (2)$$



173 To differentiate the roles of neurons across various tasks, it is essential to assess neuron activations in  
 174 relation to specific goals. Intuitively, we can consider a neuron  $\mathcal{N}$  to be positively contributing to the  
 175 goal at step  $t$  when its activation  $a(\mathcal{N}, t)$  surpasses the average activation  $\bar{a}(\mathcal{N})$ , i.e.  $a(\mathcal{N}, t) > \bar{a}(\mathcal{N})$ ,  
 176 while the GM at the same step also exceeds its average, i.e.  $q(t) > \bar{q}$ . To quantify this contribution,  
 177 we accumulate a batch of results over  $T$  steps and define the over-activation rate as follows:

$$R_{over}(\mathcal{N}) = \frac{\sum_{t=1}^T \mathbb{1}_{[a(\mathcal{N}, t) > \bar{a}(\mathcal{N})] = \mathbb{1}_{[q(t) > \bar{q}]}}}{T}. \quad (3)$$

178 Here,  $\mathbb{1}_{[condition]} \in \{0, 1\}$  denotes the indicator function, which returns 1 if and only if the specified  
 179 condition is satisfied. While Eq. (3) assesses the positive correlation of neurons towards achieving  
 180 the goal, where a higher rate implies a greater significance of the neuron in producing better outcome,  
 181 however, it overlooks neurons that exhibit a negative correlation with the goal but still carry valuable  
 182 task-related knowledge. Specifically, when the activation of a neuron falls below its average activation,  
 183 the agent performs well conversely. To this end, we define a **comprehensive score**  $Score(\mathcal{N})$  for  
 184 the neuron that takes into account both positive and negative effects:

$$Score(\mathcal{N}) = \max(R_{over}(\mathcal{N}), 1 - R_{over}(\mathcal{N})). \quad (4)$$

185 Subsequently, we rank all neurons in the agent network, excluding those in the last layer, in descending  
 186 order based on their scores. The RL skill neurons are determined by selecting the neurons with the  
 187 top  $m\%$  highest scores, formally defined as follows, where  $\tau_m(\cdot)$  denotes the top- $m$  selection operator.  
 188 And the pseudo-code of the identification method is shown in Appendix D.

$$\mathcal{N}_{RL\ skill} = \tau_m(Score(\mathcal{N})) \quad (5)$$

### 189 3.3 Neuron-level Balance between Stability and Plasticity

190 Building upon the concept of RL skill neurons, we propose a novel DRL framework — **Neuron-level**  
 191 **Balance between Stability and Plasticity (NBSP)**, as shown in Figure 2. Unlike prior methods (Bai  
 192 et al., 2023; Kim et al., 2023), the framework proposed does not require complex network designs or  
 193 additional parameters. Given that RL skill neurons encode essential task-specific skills, preserving  
 194 their activation patterns is critical to maintaining knowledge from previous tasks during continual  
 195 tasks learning. However, simply freezing RL skill neurons would hinder the ability of the agent to  
 196 adapt to new tasks. To address this challenge, NBSP employs an adaptive **gradient masking**  
 197 technique. Specifically, during each update round in the continual learning process, the gradients of  
 198 RL skill neurons are selectively masked to restrict changes in their activation patterns while allowing  
 199 other neurons to adapt freely. This process is formally expressed as follows:

$$\Delta W_{:,j} = mask_j^{(l)} \cdot \Delta W_{:,j}^{(l)}, \quad (6)$$

200 where  $\Delta W_{:,j}^{(l)}$  denotes the gradient with respect to the weight  $W_{:,j}^{(l)}$  in the  $l$ -th layer of the network,  
 201 and  $j$  is the index of the output neuron in that layer. The term  $mask_j^{(l)}$  is associated with the score of  
 202  $j$ -th neuron in the  $l$ -th layer, which could be calculated as follows:

$$mask(\mathcal{N}) = \begin{cases} \alpha(1 - Score(\mathcal{N})) & \text{if } \mathcal{N} \in \mathcal{N}_{RL\ skill} \\ 1 & \text{if } \mathcal{N} \notin \mathcal{N}_{RL\ skill} \end{cases}, \quad (7)$$

203 where  $\mathcal{N}_{RL\ skill}$  represents the set of RL skill neurons, and  $\alpha$  is a super-parameter that determines the  
 204 degree of restriction on these neurons, which is configured to 0.2 in the experiment. **By employing**  
 205 **gradient masking, NBSP effectively safeguards the encoded skills within RL skill neurons from**  
 206 **interference during the learning of new tasks, thereby enhancing stability. At the same time, RL**  
 207 **skill neurons remain adaptable, allowing fine-tuning to accommodate new tasks and maintaining**  
 208 **high plasticity. In addition, neurons except RL skill neurons are free to fully engage in learning**  
 209 **new task-specific knowledge, ensuring comprehensive learning across tasks.**

210 To mitigate excessive drift from knowledge acquired in previous tasks, we integrate the **experience**  
 211 **replay** technique, periodically sampling prior experiences at specific intervals  $k$ . After training on a  
 212 task, a portion of the experiences, rather than the entirety, are stored in a unified replay buffer  $D_{pre}$ ,  
 213 requiring only a modest memory footprint. By incorporating experience replay, the stability of DRL  
 214 agents is further enhanced. The corresponding loss function is defined as follows:

$$\mathcal{L} = R(t) \cdot \mathbb{E}_{(s_t, a_t, s_{t+1}, r_t, d_t) \sim D_{pre}} [L] + (1 - R(t)) \cdot \mathbb{E}_{(s_t, a_t, s_{t+1}, r_t, d_t) \sim D} [L], \quad (8)$$

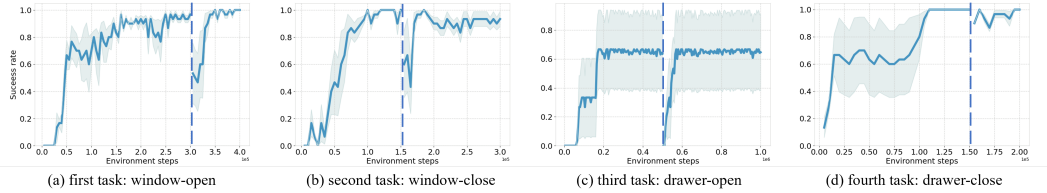


Figure 3: Training process of NBSP on the Meta-World benchmark. The segments to the left and right of the dashed line represent the training processes of the first and second cycles, respectively.

215 where  $L$  denotes the original loss function,  $R(t)$  is a binary function that evaluates to 1 if and  
 216 only if the current step  $t$  is at an interval.  $D$  represents the replay buffer for the current task, and  
 217  $(s_t, a_t, s_{t+1}, r_t, d_t)$  denotes the tuple of the current state, action, next state, reward, and whether the  
 218 episode is done sampled from the replay buffer. The pseudo-code of NBSP is shown in Appendix D.

## 219 4 Experiment

220 In this section, we evaluate the performance of NBSP on the **Meta-World** (Yu et al., 2020) and **Atari**  
 221 benchmarks (Mnih et al., 2013).

Table 1: Results of NBSP with other baselines on the Meta-World benchmark.

Cycling sequential tasks	Metrics	Methods							
		EWC	NPC	ANCL	CoTASP	CReLu	CBP	PI	NBSP
(window-open → window-close)	ASR ↑	0.63 ± 0.03	0.26 ± 0.01	0.66 ± 0.04	0.05 ± 0.01	0.26 ± 0.14	0.67 ± 0.05	0.61 ± 0.02	<b>0.90 ± 0.04</b>
	FM ↓	0.89 ± 0.07	0.68 ± 0.04	0.84 ± 0.10	<b>0.01 ± 0.01</b>	0.66 ± 0.42	0.78 ± 0.13	0.91 ± 0.07	<b>0.18 ± 0.01</b>
	FWT ↑	<b>0.97 ± 0.02</b>	0.26 ± 0.01	0.97 ± 0.03	0.04 ± 0.01	0.33 ± 0.19	0.95 ± 0.02	0.95 ± 0.01	<b>0.96 ± 0.02</b>
(drawer-open → drawer-close)	ASR ↑	0.68 ± 0.06	0.35 ± 0.05	0.64 ± 0.02	0.07 ± 0.01	0.29 ± 0.20	0.61 ± 0.03	0.60 ± 0.07	<b>0.96 ± 0.02</b>
	FM ↓	0.80 ± 0.15	0.69 ± 0.05	0.88 ± 0.09	<b>0.01 ± 0.01</b>	0.31 ± 0.32	0.91 ± 0.03	0.71 ± 0.30	<b>0.07 ± 0.06</b>
	FWT ↑	<b>0.98 ± 0.01</b>	0.39 ± 0.09	0.96 ± 0.01	0.09 ± 0.00	0.42 ± 0.28	0.93 ± 0.04	0.88 ± 0.15	<b>0.98 ± 0.01</b>
(button-press-topdown → window-open)	ASR ↑	0.66 ± 0.06	0.25 ± 0.00	0.61 ± 0.01	0.03 ± 0.00	0.33 ± 0.10	0.62 ± 0.01	0.63 ± 0.02	<b>0.95 ± 0.05</b>
	FM ↓	0.85 ± 0.14	0.67 ± 0.00	0.95 ± 0.05	<b>0.01 ± 0.00</b>	0.94 ± 0.01	0.97 ± 0.03	0.97 ± 0.05	<b>0.08 ± 0.12</b>
	FWT ↑	0.96 ± 0.01	0.25 ± 0.01	0.95 ± 0.03	0.04 ± 0.01	0.42 ± 0.20	<b>0.98 ± 0.02</b>	<b>0.98 ± 0.02</b>	<b>0.98 ± 0.01</b>
(window-open → window-close → drawer-open → drawer-close)	ASR ↑	0.44 ± 0.05	0.19 ± 0.04	0.48 ± 0.04	0.04 ± 0.01	0.10 ± 0.06	0.43 ± 0.03	0.41 ± 0.06	<b>0.66 ± 0.14</b>
	FM ↓	0.74 ± 0.11	0.50 ± 0.02	0.80 ± 0.04	<b>0.04 ± 0.01</b>	0.39 ± 0.02	0.91 ± 0.05	0.84 ± 0.05	<b>0.48 ± 0.18</b>
	FWT ↑	0.83 ± 0.10	0.20 ± 0.05	0.89 ± 0.06	0.08 ± 0.01	0.13 ± 0.10	<b>0.97 ± 0.02</b>	0.82 ± 0.10	<b>0.89 ± 0.12</b>
(button-press-topdown → window-close → door-open → drawer-close)	ASR ↑	0.43 ± 0.03	0.17 ± 0.01	0.44 ± 0.03	0.04 ± 0.01	0.14 ± 0.11	0.41 ± 0.02	0.38 ± 0.01	<b>0.74 ± 0.07</b>
	FM ↓	0.81 ± 0.09	0.47 ± 0.01	0.87 ± 0.02	<b>0.04 ± 0.00</b>	0.62 ± 0.16	0.94 ± 0.02	0.97 ± 0.02	<b>0.34 ± 0.15</b>
	FWT ↑	0.88 ± 0.10	0.19 ± 0.02	0.91 ± 0.08	0.07 ± 0.02	0.17 ± 0.15	<b>0.97 ± 0.01</b>	0.92 ± 0.07	<b>0.95 ± 0.06</b>

222 **Experiment setting.** We follow the the experimental paradigm of Abbas et al. (2023); Liu et al.  
 223 (2024), evaluating our proposed method on a **cycling sequence of tasks** characterized by non-  
 224 stationarity due to changing environments over time. Specifically, the agent learns each task sequen-  
 225 tially and transitions to the next without resetting the learned networks. The task cycles through a  
 226 fixed sequence, with a cycle completing once all tasks in the sequence have been learned. The agent  
 227 cycles twice, resulting in each task being repeated twice during the training process. Compared to  
 228 the CRL training paradigm, our cycling training paradigm provides a more specific evaluation of the  
 229 balance between stability and plasticity. By repeating each task twice within a cycling sequence, the  
 230 setup not only assesses the plasticity in adapting to new tasks but also evaluates its stability when  
 231 revisiting previously learned tasks, avoiding the influence of task order. Details about the benchmarks  
 232 are shown in Appendix C.2.

233 For all experiments, we use the Soft Actor-Critic (SAC) (Haarnoja et al., 2018) algorithm, as  
 234 implemented by CleanRL (Huang et al., 2022). Each agent is trained until either reaching a predefined  
 235 maximum number of steps or demonstrating stable mastery of the task in the Meta-World benchmark.  
 236 Each experiment is repeated using three different random seeds. The shaded regions in the figures  
 237 and the plus/minus numbers represent the standard error across multiple seeds. Detailed descriptions  
 238 of the hyper-parameters and other experimental settings are provided in Appendix C.3.

239 **Metric.** Overall performance is commonly assessed using the **Average Success Rate (ASR)**,  
 240 analogous to the AIA metric (Wang et al., 2024). Let  $sr_{i,j}$  represent the success rate on the  $j$ -th task

241 after completing the learning of the  $i$ -th task ( $i \geq j$ ),  $H$  denote the number of tasks. The ASR is  
 242 defined as follows. The higher the ASR, the better the method balances stability and plasticity.

$$ASR = \frac{1}{H} \sum_{i=1}^H \frac{1}{i} \sum_{i \geq j} sr_{i,j}, \quad (9)$$

243 To evaluate the stability of the agent, we utilize the **Forgetting Measure (FM)** (Chaudhry et al.,  
 244 2018a). The lower the FM, the better the method maintains stability, which is calculated as:

$$FM = \frac{1}{H-1} \sum_{i=2}^H \frac{1}{i-1} \sum_{i \geq j} \max_{l \in \{1, \dots, i-1\}} (sr_{l,j} - sr_{i,j}). \quad (10)$$

245 To assess the plasticity of the agent, we employ the **Forward Transfer (FWT)** metric (Lopez-Paz &  
 246 Ranzato, 2017), which is calculated as follows:

$$FWT = \frac{1}{H} \sum_{i=1}^H sr_{i,i}. \quad (11)$$

247 The higher the FWT, the better the method maintains plasticity. Further details about evaluation  
 248 metrics are available in Appendix C.4.

249 **Baseline.** To assess the effectiveness of our proposed NBSP framework, we compare it with seven  
 250 baseline methods dealing with the balance between stability and plasticity. **EWC** (Kirkpatrick et al.,  
 251 2017) and **NPC** (Paik et al., 2019) primarily emphasize maintaining stability, while **CRelu** (Abbas  
 252 et al., 2023), **CBP** (Dohare et al., 2024), and **PI** (Nikishin et al., 2024) focus on enhancing plasticity.  
 253 **ANCL** (Kim et al., 2023) and **CoTASP** (Yang et al., 2023) aim to achieve a balance between stability  
 254 and plasticity. Notably, CoTASP makes relevant tasks share more neurons in the meta-policy network,  
 255 and NPC estimates the importance value of each neuron and consolidates important neurons, they are  
 256 both relevant to neurons. Detailed descriptions of these baselines can be found in Appendix C.1.

## 257 4.1 Experiment on the Meta-World Benchmark

258 The experimental results of NBSP compared with other baselines on the Meta-World benchmark  
 259 are presented in Table 1. As shown in the final column, NBSP significantly outperforms all other  
 260 methods in the overall performance metric ASR. For two-task cycling tasks, NBSP achieves an ASR  
 261 consistently above 0.9, which is substantially higher than other baselines. Its stability metric, FM, is  
 262 markedly lower, while its plasticity metric, FWT, remains at a high level. Furthermore, NBSP also  
 263 demonstrates excellent performance in four-task cycling tasks, maintaining a substantial lead.

264 For stability-focused baselines, EWC achieves a relatively good ASR compared to other baselines  
 265 but still falls short of NBSP. Moreover, EWC exhibits poor stability due to its high FM values. NPC  
 266 performs even worse, failing to maintain both stability and plasticity effectively. Among plasticity-  
 267 focused baselines, CBP and PI achieve comparable plasticity to NBSP, as reflected in their high FWT  
 268 scores. However, both suffer from severe stability loss, indicated by their higher FM values. Another  
 269 plasticity-focused method, CRelu, underperforms in both stability and plasticity. For baselines  
 270 attempting to balance stability and plasticity, ANCL achieves high plasticity with competitive FWT  
 271 scores but fails to retain prior knowledge, as reflected by its high FM value. CoTASP, despite being  
 272 explicitly designed for this trade-off, performs poorly overall. Its low FM is attributed to a failure to  
 273 acquire meaningful task knowledge, as evidenced by its low FWT value.

274 The effectiveness of NBSP is further demonstrated in Figure 3, which showcases the training dynamics  
 275 of NBSP. Specifically, during the second cycle of learning the same task, the agent exhibits a high  
 276 success rate even before retraining, indicating that it has retained significant task knowledge. As  
 277 a result, the agent is able to master the task more rapidly. This highlights the ability of NBSP to  
 278 preserve knowledge from prior tasks while simultaneously maintaining the plasticity required to learn  
 279 new tasks effectively. The other training process is demonstrated in Appendix C.7. In summary,  
 280 ***NBSP delivers a remarkable improvement in maintaining stability without compromising plasticity,***  
 281 ***achieving a well-balanced trade-off in DRL.***

## 282 4.2 Ablation Study

283 In the ablation study, we further evaluate the effectiveness of (1) the two primary components of  
 284 NBSP: the gradient masking technique and experience replay technique, (2) the neuron identification

Table 2: Results of ablation study of gradient masking and experience replay techniques.

Metrics	(button-press-topdown $\rightarrow$ window-open)				
	vanilla SAC	only experience replay	only gradient masking	NBSP with hard gradient masking	NBSP
ASR $\uparrow$	0.62 $\pm$ 0.01	0.70 $\pm$ 0.08	0.71 $\pm$ 0.06	0.71 $\pm$ 0.03	<b>0.95 <math>\pm</math> 0.05</b>
FM $\downarrow$	0.99 $\pm$ 0.02	0.50 $\pm$ 0.16	0.73 $\pm$ 0.21	0.72 $\pm$ 0.04	<b>0.08 <math>\pm</math> 0.12</b>
FWT $\uparrow$	<b>0.98 <math>\pm</math> 0.02</b>	0.92 $\pm$ 0.05	0.97 $\pm$ 0.02	<b>0.98 <math>\pm</math> 0.03</b>	<b>0.98 <math>\pm</math> 0.01</b>

Table 3: Results of ablation study of neuron identification methods.

Metrics	(window-open $\rightarrow$ window-close)				(drawer-open $\rightarrow$ drawer-close)				(button-press-topdown $\rightarrow$ window-open)			
	activation	weight	random	ours	activation	weight	random	ours	activation	weight	random	ours
ASR $\uparrow$	0.65 $\pm$ 0.30	0.73 $\pm$ 0.20	0.78 $\pm$ 0.09	<b>0.90<math>\pm</math>0.04</b>	0.82 $\pm$ 0.06	0.51 $\pm$ 0.17	0.72 $\pm$ 0.26	<b>0.96<math>\pm</math>0.02</b>	0.75 $\pm$ 0.01	0.93 $\pm$ 0.06	0.72 $\pm$ 0.01	<b>0.95<math>\pm</math>0.05</b>
FM $\downarrow$	0.56 $\pm$ 0.37	0.44 $\pm$ 0.31	0.42 $\pm$ 0.13	<b>0.18<math>\pm</math>0.01</b>	0.44 $\pm$ 0.16	0.67 $\pm$ 0.00	0.41 $\pm$ 0.28	<b>0.07<math>\pm</math>0.06</b>	0.65 $\pm$ 0.02	0.15 $\pm$ 0.12	0.70 $\pm$ 0.05	<b>0.08<math>\pm</math>0.12</b>
FWT $\uparrow$	0.73 $\pm$ 0.35	0.81 $\pm$ 0.22	0.90 $\pm$ 0.06	<b>0.96<math>\pm</math>0.02</b>	0.98 $\pm$ 0.02	0.69 $\pm$ 0.22	0.83 $\pm$ 0.23	<b>0.98<math>\pm</math>0.01</b>	<b>0.99<math>\pm</math>0.00</b>	0.98 $\pm$ 0.02	0.96 $\pm$ 0.02	0.98 $\pm$ 0.01

285 method, and (3) the two critical modules of DRL: the actor and the critic. What’s more, we analyze  
 286 how the proportion of RL skill neurons influences the performance of NBSP.

287 **Gradient masking and experience replay.** To evaluate the contributions of the two core components  
 288 of NBSP, we designed five experimental settings: (1) vanilla SAC, (2) SAC with only the experience  
 289 replay, (3) SAC with only the gradient masking, (4) SAC with experience replay and hard gradient  
 290 masking, where the masks of RL skill neurons are set directly to zero, and (5) NBSP.

291 The results of the cycling sequential tasks (button-press-topdown  $\rightarrow$  window-open) are shown in  
 292 Table 2. From the results, we observe the following: (1) The vanilla SAC algorithm suffers from  
 293 severe stability loss, as indicated by a high FM score, underscoring the need for mechanisms to retain  
 294 prior knowledge. (2) Using either experience replay or gradient masking alone alleviates the stability  
 295 loss to some extent, confirming their individual effectiveness. (3) Combining both techniques in  
 296 NBSP significantly improves performance, with lower FM (indicating enhanced stability) and higher  
 297 FWT (demonstrating maintained plasticity). (4) Our adaptive gradient masking, which sets masks of  
 298 RL skill neurons based on their scores, outperforms hard masking (setting masks to zero directly),  
 299 demonstrating its superior effectiveness. *These findings demonstrate that neither experience replay  
 300 nor gradient masking alone can properly balance stability and plasticity, while their combination  
 301 achieves optimal performance.* The reason is that gradient masking and experience replay focus on  
 302 different mechanisms and therefore complement each other. Gradient masking primarily targets RL  
 303 skill neurons to reduce interference with past knowledge while maintaining the ability to fine-tune for  
 304 new tasks. And experience replay mainly acts on neurons except RL skill neurons to prevent these  
 305 neurons from being overly biased toward new tasks. Additional results for different task settings are  
 306 provided in Appendix C.8.

307 **Neuron identification method.** To evaluate the proposed goal-oriented neuron identification method,  
 308 we compare it with three alternative strategies: (1) random neuron identification, (2) identifying  
 309 neurons with activation magnitude (Jung et al., 2020), and (3) identifying neurons with weight  
 310 magnitude (Dohare et al., 2021). As shown in Table 3, our goal-oriented method consistently  
 311 outperforms the other three methods across all three metrics: ASR, FM, and FWT, which confirms  
 312 that our method effectively identifies neurons critical for knowledge retention, ensuring better stability  
 313 and plasticity in cycling sequential task learning. *These findings validate the necessity of task-  
 314 specific, goal-oriented neuron identification in enhancing balance between stability and plasticity.*

315 **Actor and critic.** To get a deeper understanding of the individual roles of the actor and critic in  
 316 DRL agents, we compare the performance of NBSP with that only applied on actor and critic. The  
 317 result is shown in Table 4. *The results indicate that both the actor and critic networks are essential  
 318 for striking an optimal balance between stability and plasticity. Notably, the critic proves to be  
 319 the more critical module in balancing this trade-off,* which aligns with the insight from Ma et al.  
 320 (2024) that plasticity loss in the critic serves as the principal bottleneck impeding efficient training in  
 321 DRL. *We further investigate this phenomenon by dissecting the inherent training mechanisms  
 322 of actor-critic RL methods, and draw the following key observations:* (1) Updates to the actor are  
 323 guided by feedback from the critic. Consequently, even if the RL skill neurons in the actor are masked,  
 324 they remain influenced by the critic, which may gradually adapt to the new task at the expense of  
 325 retaining prior knowledge; (2) In contrast, applying NBSP to the critic network indirectly constrains  
 326 the actor as well; and (3) The update process of the critic network is recursive, with its target network  
 327 updated via an exponential moving average, enabling it to preserve knowledge from the previous task  
 328 while integrating new skills. Therefore, NBSP achieves better performance on the critic than on the

Table 4: Results of ablation study of the actor and critic modules.

Metric	(window-open → window-close)			(drawer-open → drawer-close)			(button-press-topdown → window-open)		
	actor	critic	both	actor	critic	both	actor	critic	both
ASR ↑	0.76 ± 0.10	0.79 ± 0.05	<b>0.90 ± 0.04</b>	0.79 ± 0.05	0.86 ± 0.02	<b>0.96 ± 0.02</b>	0.81 ± 0.11	0.85 ± 0.16	<b>0.95 ± 0.05</b>
FM ↓	0.58 ± 0.19	0.48 ± 0.09	<b>0.18 ± 0.01</b>	0.55 ± 0.15	0.31 ± 0.03	<b>0.07 ± 0.06</b>	0.45 ± 0.28	0.35 ± 0.38	<b>0.08 ± 0.12</b>
FWT ↑	<b>0.97 ± 0.04</b>	0.94 ± 0.05	0.96 ± 0.02	<b>0.99 ± 0.01</b>	0.96 ± 0.02	0.98 ± 0.01	0.95 ± 0.01	0.95 ± 0.03	<b>0.98 ± 0.01</b>

Table 5: Results of NBSP with other baselines on the Atari benchmark.

Cycling sequential games	Metrics	Methods							
		EWC	NPC	ANCL	CoTASP	CRelu	CBP	PI	NBSP
(Pong → Bowling)	AR ↑	0.66 ± 0.07	0.51 ± 0.02	0.42 ± 0.29	-0.05 ± 0.02	0.02 ± 0.00	-0.09 ± 0.00	0.53 ± 0.01	<b>0.87 ± 0.01</b>
	FM ↓	0.58 ± 0.20	0.51 ± 0.04	0.46 ± 0.31	0.07 ± 0.01	<b>0.01 ± 0.00</b>	0.06 ± 0.00	0.78 ± 0.02	<b>0.05 ± 0.03</b>
	FWT ↑	0.70 ± 0.02	0.35 ± 0.02	0.47 ± 0.31	-0.05 ± 0.05	0.02 ± 0.01	-0.09 ± 0.00	0.60 ± 0.00	<b>0.72 ± 0.01</b>
(BankHeist → Alien)	AR ↑	0.46 ± 0.01	0.38 ± 0.06	0.46 ± 0.01	-0.08 ± 0.05	0.08 ± 0.05	0.12 ± 0.02	0.48 ± 0.14	<b>0.57 ± 0.02</b>
	FM ↓	0.98 ± 0.02	0.46 ± 0.14	0.98 ± 0.03	<b>0.27 ± 0.04</b>	0.52 ± 0.29	0.44 ± 0.09	0.88 ± 0.27	<b>0.65 ± 0.07</b>
	FWT ↑	0.71 ± 0.02	0.37 ± 0.03	0.72 ± 0.01	-0.16 ± 0.07	0.28 ± 0.11	0.30 ± 0.05	<b>0.73 ± 0.26</b>	<b>0.72 ± 0.05</b>

actor. This demonstrates the distinct roles of the actor and critic in balancing stability and plasticity, providing valuable insights for future research in this field.

**The proportion of RL skill neurons.** To evaluate the impact of the proportion of RL skill neurons on the performance of NBSP, we experiment with various proportions on the (button-press-topdown → window-open) cycling tasks. The results, shown in Figure 4, reveal an interesting trend: *as the proportion of RL skill neurons increases, the ASR improves initially, but begins to decline after reaching a certain threshold.* Specifically, when the proportion is small, not all neurons encoding task-specific skills are identified, leading to knowledge loss stored in neurons that are not selected. On the other hand, when the proportion becomes too large, neurons that do not encode skills may be incorrectly selected as RL skill neurons, which compromises their learning capacity and causes the true RL skill neurons to adjust their activations to accommodate new tasks, ultimately reducing stability. Thus, determining the optimal proportion of RL skill neurons is crucial for achieving the best performance. Our experiments suggest that a proportion of 0.2 is ideal for balancing stability and plasticity.

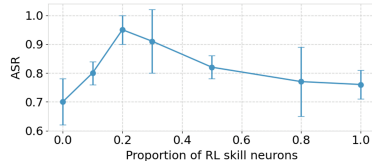


Figure 4: Performance of NBSP with different proportions of RL skill neurons.

### 4.3 Experiment on the Atari Benchmark

We further evaluate NBSP on the Atari benchmark to assess its generalization ability. In contrast to the continuous action space of Meta-World, Atari games feature discrete action spaces, and episode returns are used to evaluate the performance of each game. The results are presented in Table 5. As with the Meta-World benchmark, NBSP demonstrates superior performance in balancing stability and plasticity, outperforming other baselines across key evaluation metrics, including AR (Average Return), FM, and FWT. In a word, *NBSP exhibits excellent generalization in balance stability and plasticity across different benchmarks.*

## 5 Conclusion

This work addresses the fundamental issue of the stability-plasticity dilemma in DRL. To tackle this problem, we introduce the concept of RL skill neurons by identifying neurons that significantly contribute to knowledge retention, building upon which we then propose the Neuron-level Balance between Stability and Plasticity framework, by employing gradient masking and experience replay techniques on RL skill neurons. Experimental results on the Meta-World and Atari benchmarks demonstrate that NBSP significantly outperforms existing methods in managing the stability-plasticity trade-off. Future research could explore the application of RL skill neurons like model distillation and extend NBSP to other learning paradigms, such as supervised learning.

## 363 References

- 364 Abbas, Z., Zhao, R., Modayil, J., White, A., and Machado, M. C. Loss of plasticity in continual deep  
365 reinforcement learning. In *Conference on Lifelong Learning Agents*, pp. 620–636. PMLR, 2023.
- 366 Ahn, H., Hyeon, J., Oh, Y., Hwang, B., and Moon, T. Reset & distill: A recipe for overcoming negative transfer  
367 in continual reinforcement learning. *arXiv preprint arXiv:2403.05066*, 2024.
- 368 Anand, N. and Precup, D. Prediction and control in continual reinforcement learning. *Advances in Neural  
369 Information Processing Systems*, 36, 2024.
- 370 Andrychowicz, O. M., Baker, B., Chociej, M., Jozefowicz, R., McGrew, B., Pachocki, J., Petron, A., Plappert,  
371 M., Powell, G., Ray, A., et al. Learning dexterous in-hand manipulation. *The International Journal of  
372 Robotics Research*, 39(1):3–20, 2020.
- 373 Atkinson, C., McCane, B., Szymanski, L., and Robins, A. Pseudo-rehearsal: Achieving deep reinforcement  
374 learning without catastrophic forgetting. *Neurocomputing*, 428:291–307, 2021a.
- 375 Atkinson, C., McCane, B., Szymanski, L., and Robins, A. Pseudo-rehearsal: Achieving deep reinforcement  
376 learning without catastrophic forgetting. *Neurocomputing*, pp. 291–307, Mar 2021b. doi: 10.1016/j.neucom.  
377 2020.11.050. URL <http://dx.doi.org/10.1016/j.neucom.2020.11.050>.
- 378 Bai, F., Zhang, H., Tao, T., Wu, Z., Wang, Y., and Xu, B. Picor: Multi-task deep reinforcement learning  
379 with policy correction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp.  
380 6728–6736, 2023.
- 381 Bau, A., Belinkov, Y., Sajjad, H., Durrani, N., Dalvi, F., and Glass, J. Identifying and controlling important  
382 neurons in neural machine translation. In *International Conference on Learning Representations*, 2018.
- 383 Bau, D., Zhu, J.-Y., Strobelt, H., Lapedriza, A., Zhou, B., and Torralba, A. Understanding the role of individual  
384 units in a deep neural network. *Proceedings of the National Academy of Sciences*, 117(48):30071–30078,  
385 2020.
- 386 Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. The arcade learning environment: An evaluation  
387 platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.
- 388 Carpenter, G. A. and Grossberg, S. A massively parallel architecture for a self-organizing neural pattern  
389 recognition machine. *Computer vision, graphics, and image processing*, 37(1):54–115, 1987.
- 390 Carpenter, G. A. and Grossberg, S. Art 2: Self-organization of stable category recognition codes for analog input  
391 patterns. In *SPIE Proceedings, Intelligent Robots and Computer Vision VI*, Feb 1988. doi: 10.1117/12.942747.  
392 URL <http://dx.doi.org/10.1117/12.942747>.
- 393 Chaudhry, A., Dokania, P. K., Ajanthan, T., and Torr, P. H. Riemannian walk for incremental learning:  
394 Understanding forgetting and intransigence. In *Proceedings of the European conference on computer vision  
395 (ECCV)*, pp. 532–547, 2018a.
- 396 Chaudhry, A., Ranzato, M., Rohrbach, M., and Elhoseiny, M. Efficient lifelong learning with a-gem. In  
397 *International Conference on Learning Representations*, 2018b.
- 398 Chen, J., Wang, X., Yao, Z., Bai, Y., Hou, L., and Li, J. Finding safety neurons in large language models. *arXiv  
399 preprint arXiv:2406.14144*, 2024.
- 400 Dai, D., Dong, L., Hao, Y., Sui, Z., Chang, B., and Wei, F. Knowledge neurons in pretrained transformers. In  
401 *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long  
402 Papers)*, Jan 2022. doi: 10.18653/v1/2022.acl-long.581. URL <http://dx.doi.org/10.18653/v1/2022.acl-long.581>.
- 404 Dohare, S., Sutton, R. S., and Mahmood, A. R. Continual backprop: Stochastic gradient descent with persistent  
405 randomness. *arXiv preprint arXiv:2108.06325*, 2021.
- 406 Dohare, S., Hernandez-Garcia, J. F., Lan, Q., Rahman, P., Mahmood, A. R., and Sutton, R. S. Loss of plasticity  
407 in deep continual learning. *Nature*, 632(8026):768–774, 2024.
- 408 Dravid, A., Gandelsman, Y., Efros, A. A., and Shocher, A. Rosetta neurons: Mining the common units in a  
409 model zoo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1934–1943,  
410 2023.

- 411 eMermillod, M., eBugaiska, A., and eBONIN, P. The stability-plasticity dilemma: Investigating the continuum  
412 from catastrophic forgetting to age-limited learning effects. *Frontiers in Psychology*, *Frontiers in Psychology*,  
413 Aug 2013.
- 414 Foundation, F. Atari environments in gymnasium. [https://gymnasium.farama.org/environments/  
415 atari/](https://gymnasium.farama.org/environments/atari/), 2024. URL <https://gymnasium.farama.org/environments/atari/>. Accessed: 2024-09-  
416 14.
- 417 Goodfellow, I. J., Mirza, M., Courville, A., and Bengio, Y. An empirical investigation of catastrophic forgetting  
418 in gradient-based neural networks. *stat*, 1050:4, 2015.
- 419 Gurnee, W. and Tegmark, M. Language models represent space and time. In *The Twelfth International  
420 Conference on Learning Representations*, Oct 2023.
- 421 Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep  
422 reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–  
423 1870. PMLR, 2018.
- 424 Huang, S., Dossa, R. F. J., Ye, C., Braga, J., Chakraborty, D., Mehta, K., and AraÅšjo, J. G. Cleanrl: High-  
425 quality single-file implementations of deep reinforcement learning algorithms. *Journal of Machine Learning  
426 Research*, 23(274):1–18, 2022.
- 427 Jung, S., Ahn, H., Cha, S., and Moon, T. Continual learning with node-importance based adaptive group sparse  
428 regularization. *Advances in neural information processing systems*, 33:3647–3658, 2020.
- 429 Kim, S., Noci, L., Orvieto, A., and Hofmann, T. Achieving a better stability-plasticity trade-off via auxiliary  
430 networks in continual learning. *CVPR2023*, Mar 2023.
- 431 Kiran, B. R., Sobh, I., Talpaert, V., Mannion, P., Al Sallab, A. A., Yogamani, S., and Pérez, P. Deep reinforcement  
432 learning for autonomous driving: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 23(6):  
433 4909–4926, 2021.
- 434 Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., Milan, K., Quan, J.,  
435 Ramalho, T., Grabska-Barwinska, A., Hassabis, D., Clopath, C., Kumaran, D., and Hadsell, R. Overcoming  
436 catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, pp. 3521–3526,  
437 Mar 2017. doi: 10.1073/pnas.1611835114. URL <http://dx.doi.org/10.1073/pnas.1611835114>.
- 438 Kumar, S., Marklund, H., and Van Roy, B. Maintaining plasticity in continual learning via regenerative  
439 regularization. 2023.
- 440 Liu, J., Obando-Ceron, J., Courville, A., and Pan, L. Neuroplastic expansion in deep reinforcement learning.  
441 *arXiv preprint arXiv:2410.07994*, 2024.
- 442 Lopez-Paz, D. and Ranzato, M. Gradient episodic memory for continual learning. *Advances in neural information  
443 processing systems*, 30, 2017.
- 444 Ma, G., Li, L., Zhang, S., Liu, Z., Wang, Z., Chen, Y., Shen, L., Wang, X., and Tao, D. Revisiting plasticity in  
445 visual reinforcement learning: Data, modules and training stages. In *The Twelfth International Conference on  
446 Learning Representations*, 2024.
- 447 Mallya, A. and Lazebnik, S. Packnet: Adding multiple tasks to a single network by iterative pruning. In  
448 *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 7765–7773, 2018.
- 449 McCloskey, M. and Cohen, N. J. Catastrophic interference in connectionist networks: The sequential learning  
450 problem. In *Psychology of learning and motivation*, volume 24, pp. 109–165. Elsevier, 1989.
- 451 Mendez, J. A., van Seijen, H., and Eaton, E. Modular lifelong reinforcement learning via neural composition.  
452 *arXiv preprint arXiv:2207.00429*, 2022.
- 453 Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M. Playing  
454 atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- 455 Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M.,  
456 Fidjeland, A. K., Ostrovski, G., et al. Human-level control through deep reinforcement learning. *nature*, 518  
457 (7540):529–533, 2015.
- 458 Molchanov, P., Tyree, S., Karras, T., Aila, T., and Kautz, J. Pruning convolutional neural networks for resource  
459 efficient inference. In *International Conference on Learning Representations*, 2022.



- 460 Nikishin, E., Schwarzer, M., D’Oro, P., Bacon, P.-L., and Courville, A. The primacy bias in deep reinforcement  
461 learning. In *International conference on machine learning*, pp. 16828–16847. PMLR, 2022a.
- 462 Nikishin, E., Schwarzer, M., D’Oro, P., Bacon, P.-L., and Courville, A. The primacy bias in deep reinforcement  
463 learning. In *International conference on machine learning*, pp. 16828–16847. PMLR, 2022b.
- 464 Nikishin, E., Oh, J., Ostrovski, G., Lyle, C., Pascanu, R., Dabney, W., and Barreto, A. Deep reinforcement  
465 learning with plasticity injection. *Advances in Neural Information Processing Systems*, 36, 2024.
- 466 Paik, I., Oh, S., Kwak, T.-Y., and Kim, I. Overcoming catastrophic forgetting by neuron-level plasticity control.  
467 *AAAI2020*, Jul 2019.
- 468 Sajjad, H., Durrani, N., and Dalvi, F. Neuron-level interpretation of deep nlp models: A survey. *Transactions of*  
469 *the Association for Computational Linguistics*, 10:1285–1303, 2022.
- 470 Sokar, G., Agarwal, R., Castro, P. S., and Evci, U. The dormant neuron phenomenon in deep reinforcement  
471 learning. In *International Conference on Machine Learning*, pp. 32145–32168. PMLR, 2023.
- 472 Sutton, R. S. Reinforcement learning: An introduction. *A Bradford Book*, 2018.
- 473 Tang, T., Luo, W., Huang, H., Zhang, D., Wang, X., Zhao, X., Wei, F., and Wen, J.-R. Language-specific neurons:  
474 The key to multilingual capabilities in large language models. *arXiv preprint arXiv:2402.16438*, 2024.
- 475 Wang, L., Zhang, X., Su, H., and Zhu, J. A comprehensive survey of continual learning: theory, method and  
476 application. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- 477 Wang, X., Wen, K., Zhang, Z., Hou, L., Liu, Z., and Li, J. Finding skill neurons in pre-trained transformer-based  
478 language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language*  
479 *Processing*, pp. 11132–11152, 2022.
- 480 Wolczyk, M., Zajac, M., Pascanu, R., Kuciński, Ł., and Miłoś, P. Disentangling transfer in continual reinforce-  
481 ment learning. *Advances in Neural Information Processing Systems*, 35:6304–6317, 2022.
- 482 Xu, H., Zhan, R., Wong, D. F., and Chao, L. S. Let’s focus on neuron: Neuron-level supervised fine-tuning for  
483 large language model. *arXiv preprint arXiv:2403.11621*, 2024.
- 484 Yang, R., Xu, H., Wu, Y., and Wang, X. Multi-task reinforcement learning with soft modularization. *Advances*  
485 *in Neural Information Processing Systems*, 33:4767–4777, 2020.
- 486 Yang, Y., Zhou, T., Jiang, J., Long, G., and Shi, Y. Continual task allocation in meta-policy network via sparse  
487 prompting. In *International Conference on Machine Learning*, pp. 39623–39638. PMLR, 2023.
- 488 Yu, T., Quillen, D., He, Z., Julian, R., Hausman, K., Finn, C., and Levine, S. Meta-world: A benchmark and  
489 evaluation for multi-task and meta reinforcement learning. In *Conference on robot learning*, pp. 1094–1100.  
490 PMLR, 2020.

## 491 A Related Wrok

492 **Balance between stability and plasticity.** In DRL, the agent faces a fundamental challenge: the  
493 stability-plasticity dilemma, first introduced by Carpenter & Grossberg (1988). Recent research has  
494 proposed various strategies to address this issue by balancing stability and plasticity.

495 Replay-based methods are widely employed to enhance stability by reusing experiences from past  
496 distributions. For example, Chaudhry et al. (2018b) introduced A-GEM, which combines episodic  
497 memory to ensure that the average loss of prior tasks does not increase when learning a new task.  
498 Similarly, Wolczyk et al. (2022) proposed ClonEx-SAC, which uses actor behavioral cloning and best-  
499 return exploration to boost performance in CRL. To reduce storage requirements, pseudo-rehearsals  
500 generated from a generative model have also been proposed (Atkinson et al., 2021a).

501 Maintaining the expressiveness of neurons is key to preserving plasticity. Nikishin et al. (2022b)  
502 proposed a mechanism that periodically resets a portion of the agent’s network to counteract plasticity  
503 loss. Likewise, Nikishin et al. (2024) introduced plasticity injection, a lightweight intervention that  
504 enhances network plasticity without increasing trainable parameters or introducing prediction bias.  
505 The Reset & Distillation (R&D) framework combines resetting the online actor-critic network for new  
506 tasks with offline distillation of knowledge from previous action probabilities, effectively retaining  
507 plasticity (Ahn et al., 2024). Additionally, Abbas et al. (2023) proposed the Concatenated ReLUs  
508 (CReLUs) activation function to prevent activation collapse, thereby alleviating plasticity degradation.

509 Modularity-based approaches have shown promise in balancing stability and plasticity by decoupling  
510 task-specific and general knowledge. For instance, Anand & Precup (2024) decomposed the value  
511 function into a permanent value function, which captures persistent knowledge, and a transient  
512 value function, which facilitates rapid adaptation. Yang et al. (2020) designed a routing network to  
513 estimate task-specific routing strategies, reconfigure the base network, and combine routes using  
514 a soft modularity mechanism, making it effective for sequential tasks. Similarly, Mendez et al.  
515 (2022) proposed a compositional lifelong RL framework that uses accumulated neural components  
516 to accelerate learning for new tasks while preserving performance on past tasks via offline RL and  
517 replayed experiences.

518 **Neuron-level Research** Recent research highlights that not all neurons remain active across varying  
519 contexts, and this neuron sparsity is often positively correlated with task-specific performance (Xu  
520 et al., 2024). Building on this insight, numerous studies have focused on identifying and leveraging  
521 skill neurons to interpret network behavior and tackle specific challenges, achieving significant  
522 advancements. For example, skill neurons in pre-trained Transformers, which demonstrate strong  
523 predictive value for task labels, have been utilized for network pruning to enhance efficiency and  
524 improve transferability (Wang et al., 2022). Sokar et al. (2023) investigate dormant neurons in deep  
525 reinforcement learning and propose a method to recycle them during training. Similarly, Dravid  
526 et al. (2023) introduce Rosetta Neurons, enabling cross-class alignments and transformations without  
527 specialized training. In large language models, language-specific neurons have been identified to  
528 control output languages by selective activation or deactivation (Tang et al., 2024), while safety  
529 neurons have been analyzed to enhance safety alignment through mechanistic interpretability (Chen  
530 et al., 2024).

531 Despite these achievements, the exploration of skill neurons in DRL remains limited. Existing neuron-  
532 level approaches primarily focus on task-specific sub-network selection. For instance, CoTASP learns  
533 hierarchical dictionaries and meta-policies to generate sparse prompts and extract sub-networks  
534 as task-specific policies (Yang et al., 2023). Similarly, Mallya & Lazebnik (2018) sequentially  
535 allocate multiple tasks within a single network through iterative pruning and re-training, balancing  
536 performance and storage efficiency. Unlike these methods, our work identifies RL skill neurons  
537 specifically tailored to deep reinforcement learning, ensuring a balance between stability and plasticity  
538 by preserving the task-relevant knowledge encoded in these neurons while allowing for fine-tuning.

## 539 B Preliminary

### 540 B.1 Markov Decision Process (MDP)

541 A Markov Decision Process(MDP) is a framework used to describe a problem involving learning  
542 from actions to achieve a goal. Almost all reinforcement learning problems can be characterized

543 as a Markov Decision Process. Each MDP is defined by a tuple  $\langle S, A, P, R, \gamma \rangle$ , where  $S$  and  $A$   
544 represent state and action spaces respectively. The transition dynamics of the MDP are defined by the  
545 function  $P : S \times A \times S \rightarrow [0, 1]$ , which represents the probability of transitioning from a give state  $s$   
546 with action  $a$  to state  $s'$ . The reward function is represented by  $R : S \times A \times S \rightarrow \mathbb{R}$ , and  $\gamma \in (0, 1)$   
547 is the discount factor. At each time step  $t$ , an agent observes the state of the environment, denoted as  
548  $s_t$ , and selects an action  $a_t$  according to a policy  $\pi(a|s)$ . One time step later, the agent receives a  
549 numerical reward  $r_{t+1}$  and transitions to a new state  $s_{t+1}$ . In the simplest case, the return is the sum  
550 of the rewards when the agent–environment interaction naturally breaks into subsequences, which we  
551 refer to episodes (Sutton, 2018).

## 552 B.2 Soft Actor-Critic (SAC)

Soft Actor-Critic (SAC) is an off-policy actor-critic deep reinforcement learning algorithm that leverages maximum entropy to promote exploration. This work employs SAC to train a policy that effectively balances stability and plasticity, chosen for its sample efficiency, excellent performance, and robust stability. In this framework, the actor aims to maximize both the expected reward and the entropy of the policy. The parameters  $\phi$  of the actor are optimized by minimizing the following loss function:

$$J_{\pi}(\phi) = E_{s_t \sim D, a_t \sim \pi_{\phi}}[\alpha \log \pi_{\phi}(a_t | s_t) - Q_{\theta}(s_t, a_t)],$$

where  $D$  is the replay buffer,  $\alpha$  is the temperature parameter controlling the trade-off between exploration and exploitation,  $\theta$  denotes the parameters of the critic network,  $\pi_{\phi}$  represents the policy learned by the actor  $\phi$ , and  $Q_{\theta}$  denotes the Q-value estimated by the critic  $\theta$ . The critic network is trained to minimize the squared residual error:

$$J_Q(\theta) = E_{(s_t, a_t, s_{t+1}) \sim D} \left[ \frac{1}{2} (Q_{\theta}(s_t, a_t) - r_t - \gamma \hat{V}(s_{t+1}))^2 \right],$$

$$\hat{V}(s_t) = E_{a_t \sim \pi_{\phi}} [Q_{\theta}(s_t, a_t) - \alpha \log \pi_{\phi}(a_t | s_t)],$$

553 where  $\gamma$  represents the discount factor.

## 554 B.3 Neuron

555 In neural networks, various components, such as blocks and layers, play distinct roles. Here, we  
556 define a neuron as a single output dimension from a layer. For example, in a fully connected layer,  
557 each output dimension corresponds to a neuron. Similarly, in a convolutional layer, each output  
558 channel represents a neuron. Furthermore, following the terminology used by Sajjad et al. (2022),  
559 we classify neurons that encapsulate a single concept as focused neurons, while a group of neurons  
560 collectively representing a concept are termed group neurons.

## 561 C Experiment

### 562 C.1 Baseline

563 **EWC:** Elastic Weight Consolidation (EWC) (Kirkpatrick et al., 2017) addresses the challenge of  
564 catastrophic forgetting by allowing neural networks to retain proficiency in previously learned tasks  
565 even after a long hiatus. It achieves this by selectively slowing down learning for weights that are  
566 crucial for retaining knowledge of these tasks. This approach has demonstrated excellent performance  
567 in sequentially solving a series of classification tasks, such as those in the MNIST handwritten digit  
568 dataset, and in learning several Atari 2600 games sequentially.

569 **NPC:** Neuron-level Plasticity Control (NPC) (Paik et al., 2019) preserves the existing knowledge  
570 from the previous tasks by controlling the plasticity of the network at the neuron level. NPC estimates  
571 the importance value of each neuron and consolidates important neurons by applying lower learning  
572 rates, rather than restricting individual connection weights to stay close to the values optimized for the  
573 previous tasks. The experimental results on the several classification datasets show that neuron-level  
574 consolidation is substantially effective.

575 **ANCL:** Auxiliary Network Continual Learning (ANCL) is an innovative approach that incorporates an  
576 auxiliary network to enhance plasticity within a model that primarily emphasizes stability. Specifically,

577 this framework introduces a regularizer that effectively balances plasticity and stability, achieving  
 578 superior performance over strong baselines in both task-incremental and class-incremental learning  
 579 scenarios.

580 **CoTASP**: Continual Task Allocation via Sparse Prompting (CoTASP) (Yang et al., 2023) learns  
 581 over-complete dictionaries to produce sparse masks as prompts extracting a sub-network for each task  
 582 from a meta-policy network. Hence, relevant tasks share more neurons in the meta-policy network  
 583 due to similar prompts while cross-task interference causing forgetting is effectively restrained. It  
 584 outperforms existing continual and multi-task RL methods on all seen tasks, forgetting reduction, and  
 585 generalization to unseen tasks.

586 **CRelu**: Concatenated ReLUs (CReLUs) (Abbas et al., 2023) is a simple activation function that  
 587 concatenates the input with its negation and applies ReLU to the result. It performs effectively in  
 588 facilitating continual learning in a changing environment.

589 **CBP**: Continual BackPropagation (CBP) (Dohare et al., 2024) reinitializes a small number of units  
 590 during training, typically fewer than one per step. To prevent disruption of what the network  
 591 has already learned, only the least-used units are considered for reinitialization. It shows great  
 592 performance on Continual ImageNet and class-incremental CIFAR-100.

593 **PI**: Plasticity Injection (PI) (Nikishin et al., 2024) freeze the parameters  $\theta$  and introduce a new set  
 594 of parameters  $\theta'$  sampled from random initialization at some point in training, where the network  
 595 might have started losing plasticity. The results on Atari show that plasticity injection attains stronger  
 596 performance compared to alternative methods while being computationally efficient.

## 597 C.2 Benchmark

598 **Meta-World**. Meta-World is an open-source benchmark for meta-reinforcement learning and  
 599 multitask learning, comprising 50 distinct robotic manipulation tasks (Yu et al., 2020).

600 All tasks are executed by a simulated Sawyer robot, with the action space defined as a 2-tuple: the  
 601 change in the 3D position of the end-effector, followed by a normalized torque applied to the gripper  
 602 fingers.

603 The observation space has a consistent dimensionality of 39, although different dimensions correspond  
 604 to various aspects of each task. Typically, the observation space is represented as a 6-tuple, including  
 605 the 3D Cartesian position of the end-effector, a normalized measure of the gripper’s openness, the 3D  
 606 position and the quaternion of the first object, the 3D position and quaternion of the second object, all  
 607 previous measurements within the environment, and the 3D position of the goal.

608 The reward function for all tasks is structured and multi-component, aiding in effective policy learning  
 609 for each task component. With this design, the reward functions maintain a similar magnitudes across  
 610 tasks, generally ranging between 0 and 10. The descriptions of the six tasks used in our experiments  
 611 are listed below, and the appearance of these tasks is shown in Figure 5.

- 612 • **drawer-open**: Open a drawer, with randomized drawer positions.
- 613 • **drawer-close**: Push and close a drawer, with randomized drawer positions.
- 614 • **window-open**: Push and open a window, with randomized window positions.
- 615 • **window-close**: Push and close a window, with randomized window positions.
- 616 • **door-open**: Open a door with a revolving joint. Randomize door positions.
- 617 • **button-press-topdown**: Press a button from the top. Randomize button positions.

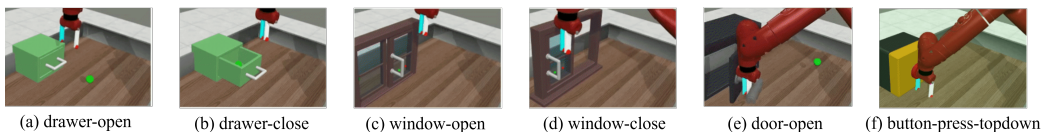


Figure 5: Tasks in the Meta-World benchmark used in our experiments.

618 **Atari.** Atari environments are simulated using the Arcade Learning Environment (ALE) (Bellemare  
619 et al., 2013) via the Stella emulator.

620 Each environment utilizes a subset of the full action space, which includes actions like NOOP,  
621 FIRE, UP, RIGHT, LEFT, DOWN, UPRIGHT, UPLEFT, DOWNRIGHT, DOWNLEFT, UPFIRE,  
622 RIGHTFIRE, LEFTFIRE, DOWNFIRE, UPRIGHTFIRE, UPLEFTFIRE, DOWNRIGHTFIRE, and  
623 DOWNLEFTFIRE. By default, most environments employ only a smaller subset of these actions,  
624 excluding those that have no effect on gameplay.

625 Observations in Atari environments are RGB images displayed to human players, with *obs\_type* =  
626 "rgb", corresponding to an observation space defined as  $Box(0, 255, (210, 160, 3), np.uint8)$ .

627 The specific reward dynamics vary depending on the environment and are typically detailed in the  
628 game's manual.

629 The descriptions of the four games used in our experiments are listed below (Foundation, 2024), and  
630 the appearance of these games is shown in Figure 6.

631 • **Bowling:** The goal is to score as many points as possible in a 10-frame game. Each frame allows  
632 up to two tries. Knocking down all pins on the first try is called a "strike", while doing so on the  
633 second try is a "spare". Failing to knock down all pins in two attempts results in an "open" frame.

634 • **Pong:** You control the right paddle and compete against the computer-controlled left paddle. The  
635 objective is to deflect the ball away from your goal and into the opponent's goal.

636 • **BankHeist:** You play as a bank robber trying to rob as many banks as possible while avoiding the  
637 police in maze-like cities. You can destroy police cars using dynamite and refill your gas tank by  
638 entering new cities. Lives are lost if you run out of gas, are caught by the police, or run over your  
639 own dynamite.

640 • **Alien:** You are trapped in a maze-like spaceship with three aliens. Your goal is to destroy their  
641 eggs scattered throughout the ship while avoiding the aliens. You have a flamethrower to fend  
642 them off and can occasionally collect a power-up (pulsar) that temporarily enables you to kill  
643 aliens.

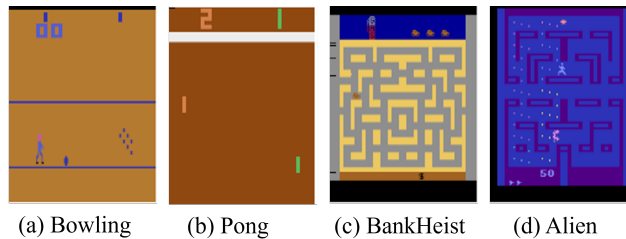


Figure 6: Games in the Atari benchmark used in our experiments.

### 644 C.3 Experiment setting

645 For all experiments, we utilize the open-source PyTorch implementation of Soft Actor-Critic (SAC)  
646 provided by CleanRL (Huang et al., 2022) on a single RTX2080Ti GPU. CleanRL is a Deep  
647 Reinforcement Learning library that offers high-quality, single-file implementations with research-  
648 friendly features. The code is both clean and straightforward, and we adhere to the configurations  
649 provided by CleanRL. During training, we employ an  $\epsilon$ -greedy exploration policy at the start,  
650 setting  $\epsilon = 1$  for the first  $10^4$  time steps to promote exploration. The environment is wrapped  
651 using Gym wrappers to facilitate experimentation. For the Meta-World benchmark, we utilize the  
652 RecordEpisodeStatistics wrapper to gather episode statistics. For the Atari benchmark, in addition  
653 to RecordEpisodeStatistics, we preprocess the  $210 \times 160$  pixel images by downsampling them to  
654  $84 \times 84$  using bilinear interpolation, converting the RGB images to the YUV format, and using only  
655 the grayscale channel. Additionally, we set a maximum limit on the number of noop and skip steps to  
656 standardize the exploration.

657 Regarding network architecture, we use the same actor and critic networks for all tasks within the  
658 same benchmark to ensure consistency. For the Meta-World benchmark, we employ a neural network

659 comprising four fully connected layers, of which the hidden size is [768, 768, 768]. For the Atari  
 660 benchmark, we use a convolutional neural network (CNN) with three convolutional layers featuring  
 661 32, 64, and 64 channels, respectively, followed by three fully connected layers, of which the hidden  
 662 size is [768, 768].

663 To reduce randomness and enhance the reliability of our results, we train each agent using three  
 664 random seeds. Additional hyper-parameters for the SAC algorithm applied in the Meta-World and  
 665 Atari benchmarks are detailed in Table 6.

Table 6: Hyper-parameters of SAC in our experiments.

Parameters	Values for Meta-World	Values for Atari
Initial collect steps	10000	20000
Discount factor	0.99	0.99
Training environment steps	$10^6$	$1.5 \times 10^6, 3 \times 10^6$
Testing environment steps	$10^5$	$10^5$
Replay buffer size	$10^6$	$2 \times 10^5$
Updates per environment step (Replay Ratio)	2	4
Target network update period	1	8000
Target smoothing coefficient	0.005	1
Optimizer	Adam	Adam
Policy learning rate	$3 \times 10^{-4}$	$10^{-4}$
Q-value learning rate	$10^{-3}$	$10^{-4}$
Minibatch size	256	64
Alpha	0.2	0.2
Autotune	True	True
Average environment steps of success rate	10	-
Stable threshold to finish training	0.9	-
Replay interval	10	10
No-op max	-	30
Target entropy scale	-	0.89
Storing experience size	$10^5$	$10^5$

#### 666 C.4 Metrics

667 For the Meta-World benchmark, the average success rate is computed over 20 episodes. For the Atari  
 668 benchmark, the success rate is replaced by the return of each episode. We normalize the return for  
 669 each game to obtain summary statistics across games, as follows:

$$R = \frac{r_{agent} - r_{random}}{r_{human} - r_{random}}, \quad (12)$$

670 where  $r_{agent}$  represents the average return evaluated over  $10^5$  steps, the random score  $r_{random}$  and  
 671 human score  $r_{human}$  are consistent with those used by Mnih et al. (2015), as detailed in Table 7.

Table 7: Normalization scores of Atari games.

games	$r_{random}$	$r_{human}$
Bowling	23.1	154.8
Pong	-20.7	9.3
BankHeist	14.2	734.4
Alien	227.5	6875

672 For the Atari benchmark tasks, the overall performance is evaluated by Average Return (AR), which  
 673 is analogous to ASR in the Meta-World benchmark. It is calculated as follows:

$$AR = \frac{1}{k} \sum_{i=1}^k \frac{1}{i} \sum_{j \geq i} R_{i,j}, \quad (13)$$

674 where  $R_{i,j}$  represents the average return evaluated on the  $j$ -th task after completing the learning of  
 675 the  $i$ -th task ( $i \geq j$ ), and  $k$  represents the number of tasks. A higher AR indicates better performance  
 676 in balancing stability and plasticity.

### 677 C.5 RL Skill Neurons

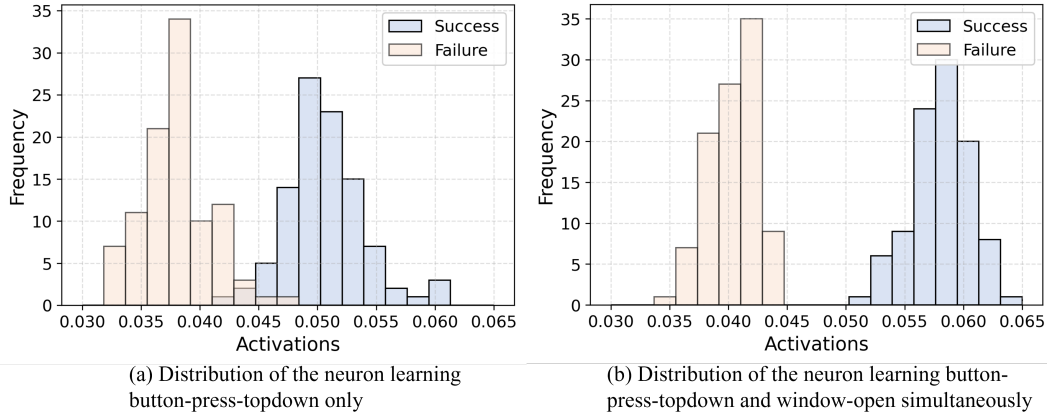


Figure 7: Distribution histogram of the activations of a neuron in two learning settings.

678 To validate the existence of RL skill neurons in sequential task learning instead of single task  
 679 learning, we conduct an additional analysis comparing the activation distributions of neurons when  
 680 learning button-press-topdown in isolation versus learning button-press-topdown and window-open  
 681 simultaneously. As shown in Figure 7, the activation distribution of a representative neuron remains  
 682 highly correlated with task success, regardless of whether it is learned in isolation or alongside  
 683 another skill. This observation supports our hypothesis that skill-specific neurons retain their essential  
 684 role even in a sequential task learning scenario.

685 Additionally, we dig deeper into the identified RL skill neurons and separate them into general and  
 686 specific skills. How to deeply investigate general skills is key for our future research. To explore  
 687 this, we design an experiment to verify the existence of general and specific skills. After sequentially  
 688 training on the button-press-topdown and window-open tasks, we identify the RL skill neurons  
 689 associated with each task. We hypothesize that the intersection set represents general skill neurons,  
 690 while the difference set represents specific skill neurons. To validate this hypothesis, we zero out the  
 691 outputs of these neurons separately. The results in Table 8 show that when the outputs of the general  
 692 skill neurons are zeroed out, the agent fails to complete both tasks. In contrast, when the outputs of  
 693 task-specific neurons are zeroed out, the agent can't complete the corresponding task but is still able  
 694 to complete the other task. This confirms the existence of both general and specific skills.

Table 8: Results of zeroing out the output of general of specific skill neurons.

tasks	zero out the in- intersection set	zero out the difference set of button-press-topdown relative to window-open	zero out the difference set of window-open relative to button-press-topdown
button-press-topdown	0	0.33	1.00
window-open	0	1.0	0.42

### 695 C.6 Results of Vanilla SAC

696 To validate the effectiveness of NBSP, it is essential to first confirm whether the vanilla SAC algorithm  
 697 can successfully solve each task individually. So we conducted experiments by training a vanilla  
 698 SAC agent on all tasks in our experiment. The results, presented in Figure 8, demonstrate that the  
 699 vanilla SAC algorithm successfully learns all tasks in our experiment. This confirms that the balance  
 700 between stability and plasticity is not an artifact of modifications to the SAC algorithm itself but



701 rather a result of NBSF. Furthermore, the failure of other methods is not due to limitations of the SAC  
 702 algorithm.

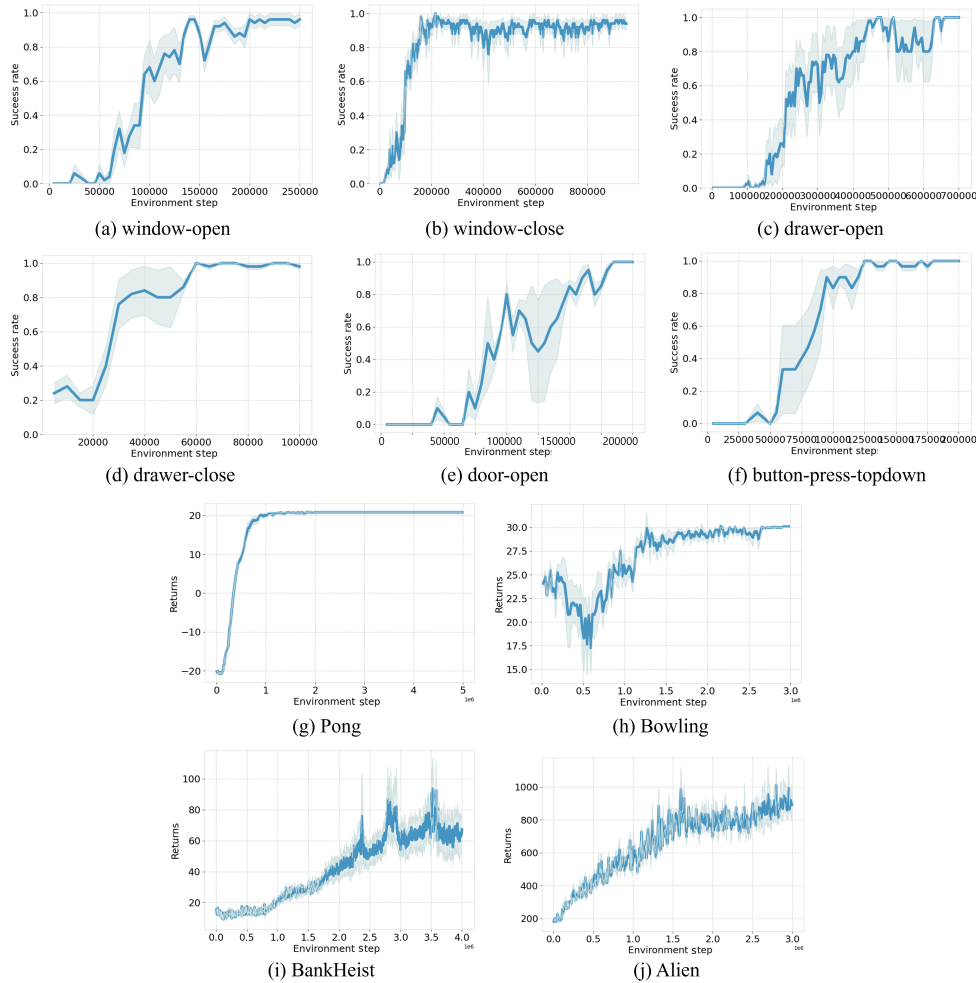


Figure 8: Training process of vanilla SAC on each individual task in our experiment.

### 703 C.7 Results on the Meta-world Benchmark

704 The training process of the other four-tasks cycling task is shown in Figure 9, and those of the  
 705 two-task cycling tasks are shown in Figure 10, Figure 11 and Figure 12 respectively. The same as  
 706 found in Section 4.1, during the second cycle of learning the same task, the agent is able to master  
 707 the task more rapidly.

### 708 C.8 Ablation Study

709 The results of the ablation study on two critical components, gradient masking and experience replay  
 710 techniques, are shown in Table 9 for the (window-open  $\rightarrow$  window-close) cycling task and in Table 10  
 711 for the (drawer-open  $\rightarrow$  drawer-close) cycling task. From these results, it is evident that both gradient  
 712 masking and experience replay techniques independently contribute to improving the stability of  
 713 the agent while maintain great plasticity. Furthermore, combining both techniques yields superior  
 714 performance, demonstrating the enhanced effectiveness of their integration.

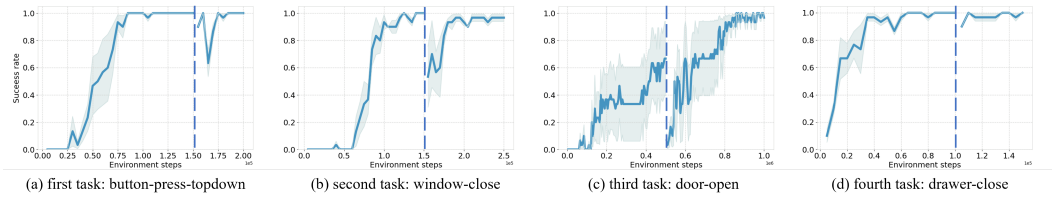


Figure 9: Training process of NBSP on (button-press-topdown  $\rightarrow$  window-close  $\rightarrow$  door-open  $\rightarrow$  drawer-close) cycling task.

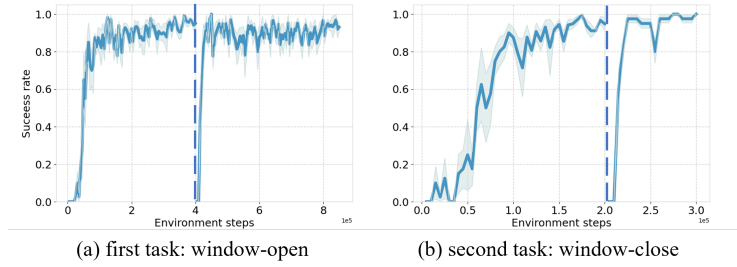


Figure 10: Training process of NBSP on (window-open  $\rightarrow$  window-close) cycling task.

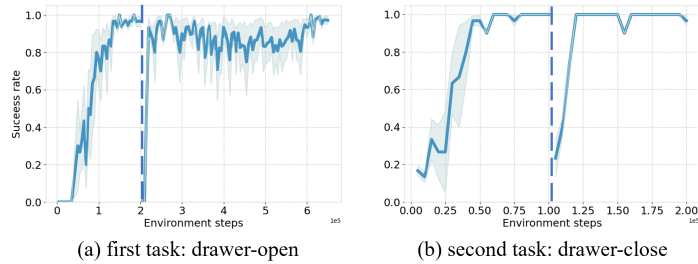


Figure 11: Training process of NBSP on (drawer-open  $\rightarrow$  drawer-close) cycling task.

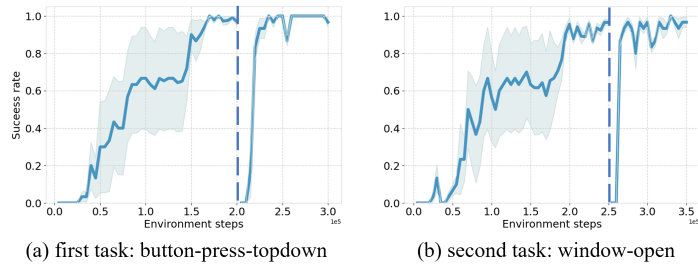


Figure 12: Training process of NBSP on (button-press-topdown  $\rightarrow$  window-open) cycling task.

Table 9: Results of ablation study of gradient masking and experience replay techniques on (window-open  $\rightarrow$  window-close) cycling task.

Metrics	(button-press-topdown $\rightarrow$ window-open)				
	vanilla SAC	only experience replay	only gradient masking	NBSP with hard gradient masking	NBSP
ASR $\uparrow$	0.63 $\pm$ 0.02	0.81 $\pm$ 0.08	0.78 $\pm$ 0.11	0.71 $\pm$ 0.04	<b>0.90 <math>\pm</math> 0.04</b>
FM $\downarrow$	0.91 $\pm$ 0.10	0.41 $\pm$ 0.13	0.54 $\pm$ 0.26	0.54 $\pm$ 0.13	<b>0.18 <math>\pm</math> 0.01</b>
FWT $\uparrow$	0.97 $\pm$ 0.02	0.96 $\pm$ 0.01	<b>0.98 <math>\pm</math> 0.01</b>	0.91 $\pm$ 0.05	<b>0.96 <math>\pm</math> 0.02</b>

Table 10: Results of ablation study of gradient masking and experience replay techniques on (drawer-open  $\rightarrow$  drawer-close) cycling task.

Metrics	(button-press-topdown $\rightarrow$ window-open)				
	vanilla SAC	only experience replay	only gradient masking	NBSP with hard gradient masking	NBSP
ASR $\uparrow$	0.67 $\pm$ 0.05	0.78 $\pm$ 0.04	0.74 $\pm$ 0.01	0.59 $\pm$ 0.16	<b>0.96 <math>\pm</math> 0.02</b>
FM $\downarrow$	0.78 $\pm$ 0.10	0.48 $\pm$ 0.10	0.64 $\pm$ 0.01	0.52 $\pm$ 0.35	<b>0.07 <math>\pm</math> 0.06</b>
FWT $\uparrow$	0.94 $\pm$ 0.04	0.97 $\pm$ 0.01	<b>0.98 <math>\pm</math> 0.02</b>	0.82 $\pm$ 0.21	<b>0.98 <math>\pm</math> 0.01</b>

## 715 D Algorithm

716 The pseudo-code of the goal-oriented method to find RL skill neurons is presented in Algorithm  
 717 1. And the pseudo-code for SAC with NBSP is presented in Algorithm 2. Key differences from  
 718 standard SAC are highlighted in blue. In addition to the extra input, two main modifications include  
 719 the sampling process and the network update process.

## 720 E Limitation and Future Work

721 **Limitation.** While the proposed NBSP method effectively balances stability and plasticity in DRL,  
 722 it does have a notable limitation. Specifically, the number of RL skill neurons must be manually  
 723 determined and adjusted according to the complexity of the learning task, as there is no automatic  
 724 mechanism for this selection. And our method currently faces challenges when applied to longer  
 725 task sequences (e.g., 10+ tasks). One key limitation is the constraint imposed by the model scale,  
 726 which inherently limits the number of skills it can learn. As the number of tasks increases, the overlap  
 727 between skill neurons across different tasks may become significant. Consequently, applying a mask  
 728 to protect RL skill neurons can restrict the learning of new tasks, making it difficult to scale without  
 729 introducing interference with previously learned knowledge.

730 **Future work.** The neuron analysis introduced in this work offers a novel approach for identifying  
 731 RL skill neurons, significantly enhancing the balance between stability and plasticity in DRL. The  
 732 identification of RL skill neurons opens up several promising directions for future research and  
 733 applications, such as: (1) Model Distillation: by focusing on RL skill neurons, it becomes possible to  
 734 distill models by pruning less relevant neurons, leading to more efficient and compact models with  
 735 minimal performance degradation. (2) Bias Control and Model Manipulation: RL skill neurons could  
 736 be leveraged to control biases and modify model behaviors by selectively adjusting their activations.  
 737 This approach could be particularly valuable in scenarios requiring specific outputs or behaviors.

738 While our current method may not yet fully address longer task sequences, it lays a strong foundation  
 739 for future research. Moving forward, we aim to explore strategies to better leverage RL skill neurons  
 740 for continual learning over an extended sequence of tasks. What’s more, its applicable potential  
 741 extends beyond DRL. It could also be adapted to other learning paradigms, such as supervised and  
 742 unsupervised learning, to address similar stability-plasticity challenges. In future work, we plan to  
 743 explore these extensions and verify their effectiveness across various domains.

---

**Algorithm 1** Procedure for Identifying RL Skill Neurons

---

**Input:** Initial average step  $T_{avg}$ , initial evaluation step  $T$ , initial proportion of RL skill neuron  $m$ , initial average activation  $\bar{a}(\mathcal{N}) = 0$ , initial average GM  $\bar{q} = 0$ , initial over-activation rate  $R_{over} = 0$ .

- 1: **for** each step  $t$  **do**
- 2:   Compute activation  $a(\mathcal{N}, t) \leftarrow \phi(\cdot)$
- 3:   Compute GM  $q(t)$
- 4:   Compute average activation:

$$\bar{a}(\mathcal{N}) = \bar{a}(\mathcal{N}) + \frac{1}{T_{avg}}a(\mathcal{N}, t).$$

- 5:   Compute average GM:

$$\bar{q} = \bar{q} + \frac{1}{T_{avg}}q(t).$$

- 6: **end for**
- 7: **for** each step  $t$  **do**
- 8:   Compute activation  $a(\mathcal{N}, t) \leftarrow \phi(\cdot)$
- 9:   Compute GM  $q(t)$
- 10:   Capture association:

$$R_{over} = R_{over} + \frac{1}{T}1_{[1_{[a(\mathcal{N}, t) > \bar{a}(\mathcal{N})]} = 1_{[q(t) > \bar{q}]}]}$$

- 11: **end for**
- 12: Derive scores  $Score$  for each neuron:

$$Score(\mathcal{N}) = \max(R_{over}(\mathcal{N}), 1 - R_{over}(\mathcal{N}))$$

- 13: Identify the top-performing neurons as RL skill neurons:

$$\mathcal{N}_{RL\ skill} = \tau_m(Score(\mathcal{N}))$$

---

---

**Algorithm 2** Neuron-level Balance between Stability and Plasticity (NBSP) Applied in SAC

---

Initialize policy parameters  $\theta$ , Q-function parameters  $\phi_1, \phi_2$ , and target Q-function parameters  $\phi'_1, \phi'_2$ Initialize empty replay buffer  $\mathcal{D}$ Initialize replay interval  $k$ **Input:** Replay buffer  $\mathcal{D}_{\text{pre}}$ , mask of the policy  $\text{mask}_\theta$  and mask of the Q-function parameters  $\text{mask}_{\phi_1}, \text{mask}_{\phi_2}$ 

```
1: for each task do
2:   for each iteration do
3:     for each environment step do
4:       Sample action  $a_t \sim \pi_\theta(a_t|s_t)$ 
5:       Execute action  $a_t$  and observe reward  $r_t$  and next state  $s_{t+1}$ 
6:       Store  $(s_t, a_t, r_t, s_{t+1})$  in replay buffer  $\mathcal{D}$ 
7:     end for
8:     for each gradient step do
9:       if  $\text{step} \equiv 0 \pmod{k}$  then Sample batch of transitions  $(s_i, a_i, r_i, s_{i+1})$  from  $\mathcal{D}_{\text{pre}}$ 
10:      else Sample batch of transitions  $(s_i, a_i, r_i, s_{i+1})$  from  $\mathcal{D}$ 
11:      end if
12:      Compute target value:
```

$$y_i = r_i + \gamma \left( \min_{j=1,2} Q_{\phi'_j}(s_{i+1}, \tilde{a}_{i+1}) - \alpha \log \pi_\theta(\tilde{a}_{i+1}|s_{i+1}) \right), \text{ where } \tilde{a}_{i+1} \sim \pi_\theta(\cdot|s_{i+1})$$

13: Update Q-functions by one step of gradient descent with mask:

$$\phi_j \leftarrow \phi_j - \lambda_Q \text{mask}_{\phi_j} \nabla_{\phi_j} \frac{1}{N} \sum_i (Q_{\phi_j}(s_i, a_i) - y_i)^2 \quad \text{for } j = 1, 2$$

14: Update policy by one step of gradient ascent with mask:

$$\theta \leftarrow \theta + \lambda_\pi \text{mask}_\theta \nabla_\theta \frac{1}{N} \sum_i \left( \alpha \log \pi_\theta(a_i|s_i) - \min_{j=1,2} Q_{\phi_j}(s_i, a_i) \right)$$

15: Update temperature  $\alpha$  by one step of gradient descent:

$$\alpha \leftarrow \alpha - \lambda_\alpha \nabla_\alpha \frac{1}{N} \sum_i (-\alpha \log \pi_\theta(a_i|s_i) - \alpha \bar{\mathcal{H}})$$

16: Update target Q-function parameters:

$$\phi'_j \leftarrow \tau \phi_j + (1 - \tau) \phi'_j \quad \text{for } j = 1, 2$$

17: **end for**18: **end for**19: **Select RL skill neurons**  $\{\mathcal{N}_{\text{RL skill}}\}$  **according to Algorithm 1**20: **Update**  $\text{mask}_{\phi_1}, \text{mask}_{\phi_2}$  **and**  $\text{mask}_\theta$ :

$$\text{mask}(\mathcal{N}) = \begin{cases} \alpha(1 - \text{Score}(\mathcal{N})) & \text{if } \mathcal{N} \in \mathcal{N}_{\text{RL skill}} \\ 1 & \text{if } \mathcal{N} \notin \mathcal{N}_{\text{RL skill}} \end{cases}$$

21: **Store part of**  $\mathcal{D}$  **into**  $\mathcal{D}_{\text{pre}}$ 22: **end for**

---

744 **NeurIPS Paper Checklist**

745 **1. Claims**

746 Question: Do the main claims made in the abstract and introduction accurately reflect the paper's  
747 contributions and scope?

748 Answer: [Yes]

749 Justification: Our main claims are summarized in Figure 2, Section 3 and Section 4 offer detailed  
750 explanations.

751 Guidelines:

- 752 • The answer NA means that the abstract and introduction do not include the claims made in  
753 the paper.
- 754 • The abstract and/or introduction should clearly state the claims made, including the contribu-  
755 tions made in the paper and important assumptions and limitations. A No or NA answer to  
756 this question will not be perceived well by the reviewers.
- 757 • The claims made should match theoretical and experimental results, and reflect how much  
758 the results can be expected to generalize to other settings.
- 759 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are  
760 not attained by the paper.

761 **2. Limitations**

762 Question: Does the paper discuss the limitations of the work performed by the authors?

763 Answer: [Yes]

764 Justification: We discuss the limitations of the work in Appendix E, and the problem setup of our  
765 work is described in Subsection 3.1.

766 Guidelines:

- 767 • The answer NA means that the paper has no limitation while the answer No means that the  
768 paper has limitations, but those are not discussed in the paper.
- 769 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 770 • The paper should point out any strong assumptions and how robust the results are to vi-  
771 olations of these assumptions (e.g., independence assumptions, noiseless settings, model  
772 well-specification, asymptotic approximations only holding locally). The authors should  
773 reflect on how these assumptions might be violated in practice and what the implications  
774 would be.
- 775 • The authors should reflect on the scope of the claims made, e.g., if the approach was only  
776 tested on a few datasets or with a few runs. In general, empirical results often depend on  
777 implicit assumptions, which should be articulated.
- 778 • The authors should reflect on the factors that influence the performance of the approach. For  
779 example, a facial recognition algorithm may perform poorly when image resolution is low or  
780 images are taken in low lighting. Or a speech-to-text system might not be used reliably to  
781 provide closed captions for online lectures because it fails to handle technical jargon.
- 782 • The authors should discuss the computational efficiency of the proposed algorithms and how  
783 they scale with dataset size.
- 784 • If applicable, the authors should discuss possible limitations of their approach to address  
785 problems of privacy and fairness.
- 786 • While the authors might fear that complete honesty about limitations might be used by  
787 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
788 limitations that aren't acknowledged in the paper. The authors should use their best judgment  
789 and recognize that individual actions in favor of transparency play an important role in  
790 developing norms that preserve the integrity of the community. Reviewers will be specifically  
791 instructed to not penalize honesty concerning limitations.

792 **3. Theory assumptions and proofs**

793 Question: For each theoretical result, does the paper provide the full set of assumptions and a  
794 complete (and correct) proof?

795 Answer: [NA]

796 Justification: Our paper is not a theoretical work.

797 Guidelines:

- 798 • The answer NA means that the paper does not include theoretical results.
- 799 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 800 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 801 • The proofs can either appear in the main paper or the supplemental material, but if they
- 802 appear in the supplemental material, the authors are encouraged to provide a short proof
- 803 sketch to provide intuition.
- 804 • Inversely, any informal proof provided in the core of the paper should be complemented by
- 805 formal proofs provided in appendix or supplemental material.
- 806 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 807 4. Experimental result reproducibility

808 Question: Does the paper fully disclose all the information needed to reproduce the main  
809 experimental results of the paper to the extent that it affects the main claims and/or conclusions of  
810 the paper (regardless of whether the code and data are provided or not)?

811 Answer: [Yes]

812 Justification: We carefully introduce our proposed framework in Section 3 and explained our  
813 settings and hyper-parameters in Section 4 and Appendix C.3.

814 Guidelines:

- 815 • The answer NA means that the paper does not include experiments.
- 816 • If the paper includes experiments, a No answer to this question will not be perceived well by
- 817 the reviewers: Making the paper reproducible is important, regardless of whether the code
- 818 and data are provided or not.
- 819 • If the contribution is a dataset and/or model, the authors should describe the steps taken to
- 820 make their results reproducible or verifiable.
- 821 • Depending on the contribution, reproducibility can be accomplished in various ways. For
- 822 example, if the contribution is a novel architecture, describing the architecture fully might
- 823 suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary
- 824 to either make it possible for others to replicate the model with the same dataset, or provide
- 825 access to the model. In general, releasing code and data is often one good way to accomplish
- 826 this, but reproducibility can also be provided via detailed instructions for how to replicate
- 827 the results, access to a hosted model (e.g., in the case of a large language model), releasing
- 828 of a model checkpoint, or other means that are appropriate to the research performed.
- 829 • While NeurIPS does not require releasing code, the conference does require all submissions
- 830 to provide some reasonable avenue for reproducibility, which may depend on the nature of
- 831 the contribution. For example
  - 832 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to
  - 833 reproduce that algorithm.
  - 834 (b) If the contribution is primarily a new model architecture, the paper should describe the
  - 835 architecture clearly and fully.
  - 836 (c) If the contribution is a new model (e.g., a large language model), then there should either
  - 837 be a way to access this model for reproducing the results or a way to reproduce the model
  - 838 (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - 839 (d) We recognize that reproducibility may be tricky in some cases, in which case authors are
  - 840 welcome to describe the particular way they provide for reproducibility. In the case of
  - 841 closed-source models, it may be that access to the model is limited in some way (e.g.,
  - 842 to registered users), but it should be possible for other researchers to have some path to
  - 843 reproducing or verifying the results.

#### 844 5. Open access to data and code

845 Question: Does the paper provide open access to the data and code, with sufficient instructions to  
846 faithfully reproduce the main experimental results, as described in supplemental material?

847 Answer: [Yes]

848 Justification: We provide our code in the supplemental material.



849

#### Guidelines:

- 850 • The answer NA means that paper does not include experiments requiring code.
- 851 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 852
- 853 • While we encourage the release of code and data, we understand that this might not be
- 854 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including
- 855 code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- 856 • The instructions should contain the exact command and environment needed to run to
- 857 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 858
- 859 • The authors should provide instructions on data access and preparation, including how to
- 860 access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 861 • The authors should provide scripts to reproduce all experimental results for the new proposed
- 862 method and baselines. If only a subset of experiments are reproducible, they should state
- 863 which ones are omitted from the script and why.
- 864 • At submission time, to preserve anonymity, the authors should release anonymized versions
- 865 (if applicable).
- 866 • Providing as much information as possible in supplemental material (appended to the paper)
- 867 is recommended, but including URLs to data and code is permitted.

#### 6. Experimental setting/details

869 Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters,  
870 how they were chosen, type of optimizer, etc.) necessary to understand the results?

871 Answer: [Yes]

872 Justification: The training and test details are described in Section 4 and Appendix C.3, and we  
873 provide the code in the supplemental material.

#### Guidelines:

- 875 • The answer NA means that the paper does not include experiments.
- 876 • The experimental setting should be presented in the core of the paper to a level of detail that
- 877 is necessary to appreciate the results and make sense of them.
- 878 • The full details can be provided either with the code, in appendix, or as supplemental
- 879 material.

#### 7. Experiment statistical significance

881 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
882 information about the statistical significance of the experiments?

883 Answer: [Yes]

884 Justification: We show the standard error in the training curves and the table of results with an  
885 average of over three random seeds.

#### Guidelines:

- 887 • The answer NA means that the paper does not include experiments.
- 888 • The authors should answer "Yes" if the results are accompanied by error bars, confidence
- 889 intervals, or statistical significance tests, at least for the experiments that support the main
- 890 claims of the paper.
- 891 • The factors of variability that the error bars are capturing should be clearly stated (for
- 892 example, train/test split, initialization, random drawing of some parameter, or overall run
- 893 with given experimental conditions).
- 894 • The method for calculating the error bars should be explained (closed form formula, call to a
- 895 library function, bootstrap, etc.)
- 896 • The assumptions made should be given (e.g., Normally distributed errors).
- 897 • It should be clear whether the error bar is the standard deviation or the standard error of the
- 898 mean.
- 899 • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably
- 900 report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality
- 901 of errors is not verified.

- 902 • For asymmetric distributions, the authors should be careful not to show in tables or figures  
903 symmetric error bars that would yield results that are out of range (e.g. negative error rates).  
904 • If error bars are reported in tables or plots, The authors should explain in the text how they  
905 were calculated and reference the corresponding figures or tables in the text.

## 906 8. Experiments compute resources

907 Question: For each experiment, does the paper provide sufficient information on the computer  
908 resources (type of compute workers, memory, time of execution) needed to reproduce the experi-  
909 ments?

910 Answer: [Yes]

911 Justification: We provide sufficient information on the computer resources in Appendix C.3.

912 Guidelines:

- 913 • The answer NA means that the paper does not include experiments.
- 914 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or  
915 cloud provider, including relevant memory and storage.
- 916 • The paper should provide the amount of compute required for each of the individual experi-  
917 mental runs as well as estimate the total compute.
- 918 • The paper should disclose whether the full research project required more compute than the  
919 experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it  
920 into the paper).

## 921 9. Code of ethics

922 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS  
923 Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

924 Answer: [Yes]

925 Justification: We have read and understood the code of ethics and have done our best to conform.

926 Guidelines:

- 927 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 928 • If the authors answer No, they should explain the special circumstances that require a  
929 deviation from the Code of Ethics.
- 930 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration  
931 due to laws or regulations in their jurisdiction).

## 932 10. Broader impacts

933 Question: Does the paper discuss both potential positive societal impacts and negative societal  
934 impacts of the work performed?

935 Answer: [NA]

936 Justification: Our work proposes a neuron-level framework to balance stability and plasticity in  
937 DRL, which does no impact the society at large, beyond improving our understanding of certain  
938 aspects of deep learning.

939 Guidelines:

- 940 • The answer NA means that there is no societal impact of the work performed.
- 941 • If the authors answer NA or No, they should explain why their work has no societal impact  
942 or why the paper does not address societal impact.
- 943 • Examples of negative societal impacts include potential malicious or unintended uses (e.g.,  
944 disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deploy-  
945 ment of technologies that could make decisions that unfairly impact specific groups), privacy  
946 considerations, and security considerations.
- 947 • The conference expects that many papers will be foundational research and not tied to  
948 particular applications, let alone deployments. However, if there is a direct path to any  
949 negative applications, the authors should point it out. For example, it is legitimate to point  
950 out that an improvement in the quality of generative models could be used to generate  
951 deepfakes for disinformation. On the other hand, it is not needed to point out that a generic  
952 algorithm for optimizing neural networks could enable people to train models that generate  
953 Deepfakes faster.

- 954 • The authors should consider possible harms that could arise when the technology is being  
955 used as intended and functioning correctly, harms that could arise when the technology is  
956 being used as intended but gives incorrect results, and harms following from (intentional or  
957 unintentional) misuse of the technology.
- 958 • If there are negative societal impacts, the authors could also discuss possible mitigation  
959 strategies (e.g., gated release of models, providing defenses in addition to attacks, mecha-  
960 nisms for monitoring misuse, mechanisms to monitor how a system learns from feedback  
961 over time, improving the efficiency and accessibility of ML).

## 962 11. Safeguards

963 Question: Does the paper describe safeguards that have been put in place for responsible release  
964 of data or models that have a high risk for misuse (e.g., pretrained language models, image  
965 generators, or scraped datasets)?

966 Answer: [NA]

967 Justification: Our paper poses no such risks for a novel framework to balance stability and  
968 plasticity in DRL.

969 Guidelines:

- 970 • The answer NA means that the paper poses no such risks.
- 971 • Released models that have a high risk for misuse or dual-use should be released with  
972 necessary safeguards to allow for controlled use of the model, for example by requiring that  
973 users adhere to usage guidelines or restrictions to access the model or implementing safety  
974 filters.
- 975 • Datasets that have been scraped from the Internet could pose safety risks. The authors should  
976 describe how they avoided releasing unsafe images.
- 977 • We recognize that providing effective safeguards is challenging, and many papers do not  
978 require this, but we encourage authors to take this into account and make a best faith effort.

## 979 12. Licenses for existing assets

980 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the  
981 paper, properly credited and are the license and terms of use explicitly mentioned and properly  
982 respected?

983 Answer: [Yes]

984 Justification: We describe the benchmarks in our experiments in Section 4 and provide the code  
985 base in Appendix C.3.

986 Guidelines:

- 987 • The answer NA means that the paper does not use existing assets.
- 988 • The authors should cite the original paper that produced the code package or dataset.
- 989 • The authors should state which version of the asset is used and, if possible, include a URL.
- 990 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 991 • For scraped data from a particular source (e.g., website), the copyright and terms of service  
992 of that source should be provided.
- 993 • If assets are released, the license, copyright information, and terms of use in the package  
994 should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated  
995 licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- 996 • For existing datasets that are re-packaged, both the original license and the license of the  
997 derived asset (if it has changed) should be provided.
- 998 • If this information is not available online, the authors are encouraged to reach out to the  
999 asset's creators.

## 1000 13. New assets

1001 Question: Are new assets introduced in the paper well documented and is the documentation  
1002 provided alongside the assets?

1003 Answer: [NA]

1004 Justification: We do not release new assets currently.

- 1005 Guidelines:
- 1006 • The answer NA means that the paper does not release new assets.
  - 1007 • Researchers should communicate the details of the dataset/code/model as part of their
  - 1008 submissions via structured templates. This includes details about training, license, limitations,
  - 1009 etc.
  - 1010 • The paper should discuss whether and how consent was obtained from people whose asset is
  - 1011 used.
  - 1012 • At submission time, remember to anonymize your assets (if applicable). You can either
  - 1013 create an anonymized URL or include an anonymized zip file.

1014 **14. Crowdsourcing and research with human subjects**

1015 Question: For crowdsourcing experiments and research with human subjects, does the paper

1016 include the full text of instructions given to participants and screenshots, if applicable, as well as

1017 details about compensation (if any)?

1018 Answer: [NA]

1019 Justification: This work does not involve crowdsourcing nor research with human subjects.

1020 Guidelines:

- 1021 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 1022 human subjects.
- 1023 • Including this information in the supplemental material is fine, but if the main contribution
- 1024 of the paper involves human subjects, then as much detail as possible should be included in
- 1025 the main paper.
- 1026 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or
- 1027 other labor should be paid at least the minimum wage in the country of the data collector.

1028 **15. Institutional review board (IRB) approvals or equivalent for research with human subjects**

1029 Question: Does the paper describe potential risks incurred by study participants, whether such

1030 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals

1031 (or an equivalent approval/review based on the requirements of your country or institution) were

1032 obtained?

1033 Answer: [NA]

1034 Justification: This work does not involve crowdsourcing nor research with human subjects.

1035 Guidelines:

- 1036 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 1037 human subjects.
- 1038 • Depending on the country in which research is conducted, IRB approval (or equivalent) may
- 1039 be required for any human subjects research. If you obtained IRB approval, you should
- 1040 clearly state this in the paper.
- 1041 • We recognize that the procedures for this may vary significantly between institutions and
- 1042 locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines
- 1043 for their institution.
- 1044 • For initial submissions, do not include any information that would break anonymity (if
- 1045 applicable), such as the institution conducting the review.

1046 **16. Declaration of LLM usage**

1047 Question: Does the paper describe the usage of LLMs if it is an important, original, or non-

1048 standard component of the core methods in this research? Note that if the LLM is used only for

1049 writing, editing, or formatting purposes and does not impact the core methodology, scientific

1050 rigor, or originality of the research, declaration is not required.

1051 Answer: [NA]

1052 Justification: The core method development in our work does not involve LLMs as any important,

1053 original, or non-standard components.

1054 Guidelines:

- 1055 • The answer NA means that the core method development in this research does not involve
- 1056 LLMs as any important, original, or non-standard components.

1057  
1058

- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.