# Triplets Better Than Pairs: Towards Stable and Effective Self-Play Fine-Tuning for LLMs

Yibo Wang<sup>1,2,\*</sup>, Hai-Long Sun<sup>1,2,\*</sup>, Guangda Huzhang, Qing-Guo Chen<sup>3</sup>, Zhao Xu<sup>3</sup>, Weihua Luo<sup>3</sup>, Kaifu Zhang<sup>3</sup>, Lijun Zhang<sup>1,4,2,†</sup>

<sup>1</sup>National Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, China

<sup>2</sup>School of Artificial Intelligence, Nanjing University, Nanjing, China

<sup>3</sup>Alibaba International Digital Commerce <sup>4</sup>Pazhou Laboratory (Huangpu), Guangzhou, China {wangyb, sunhl, zhanglj}@lamda.nju.edu.cn
{qingguo.cqg, changgong.xz, weihua.luowh, kaifu.zkf}@alibaba-inc.com

#### **Abstract**

Recently, self-play fine-tuning (SPIN) has been proposed to adapt large language models to downstream applications with scarce expert-annotated data, by iteratively generating synthetic responses from the model itself. However, SPIN is designed to optimize the current reward advantages of annotated responses over synthetic responses at hand, which may gradually vanish during iterations, leading to unstable optimization. Moreover, the utilization of reference policy induces a misalignment issue between the reward formulation for training and the metric for generation. To address these limitations, we propose a novel Triplet-based Self-Play fIne-tuNing (T-SPIN) method that integrates two key designs. First, beyond current advantages, T-SPIN additionally incorporates historical advantages between iteratively generated responses and proto-synthetic responses produced by the initial policy. Even if the current advantages diminish, historical advantages remain effective, stabilizing the overall optimization. Second, T-SPIN introduces the entropy constraint into the self-play framework, which is theoretically justified to support reference-free fine-tuning, eliminating the training-generation discrepancy. Empirical results on various tasks demonstrate not only the superior performance of T-SPIN over SPIN, but also its stable evolution during iterations. Remarkably, compared to supervised fine-tuning, T-SPIN achieves comparable or even better performance with only 25% samples, highlighting its effectiveness when faced with scarce annotated data.

#### 1 Introduction

Large language models (LLMs) have demonstrated powerful generalization capabilities and profound impact across various fields [Ouyang et al., 2022, Huang and Chang, 2023, Qin et al., 2023, Bai et al., 2023, Achiam et al., 2023, Liu et al., 2023, Jiang et al., 2024a, Wu et al., 2024a, Zhao et al., 2024, Chen et al., 2024a, Liu et al., 2024, Hui et al., 2024, Guo et al., 2024a, Jiang et al., 2024b, Anthropic, 2024a, Google, 2024, OpenAI, 2024, Sun et al., 2025a,b, Lu et al., 2025, Jiang et al., 2025]. The remarkable success of LLMs is inherently linked to the utilization of vast expert annotations. Nevertheless, considerable curation costs for high-quality data pose major obstacles to downstreamtask adaptation [Wang et al., 2024a, Tan et al., 2024]. Consequently, numerous efforts have been devoted to addressing this challenge, such as data augmentation [Dai et al., 2023, Zhu et al., 2023, Guo et al., 2024b] and model distillation [Gunasekar et al., 2023, Luo et al., 2024, Xin et al., 2024].

<sup>\*</sup>Work done during the internship at Alibaba International Digital Commerce.

<sup>&</sup>lt;sup>†</sup>Lijun Zhang is the corresponding author.

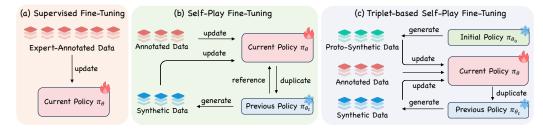


Figure 1: Comparisons of three strategies: (a) supervised fine-tuning requires large amounts of annotated data to train  $\pi_{\theta}$ ; (b) self-play fine-tuning operates with limited annotated data and iteratively generated samples, and employs the previous policy  $\pi_{\theta_t}$  as a reference for updates; (c) triplet-based self-play fine-tuning employs triplet inputs, i.e., annotated data, synthetic samples, and proto-synthetic ones from the initial policy  $\pi_{\theta_0}$ , and updates  $\pi_{\theta}$  without auxiliaries from any reference policies.

Recently, a promising research direction known as self-play fine-tuning has been introduced [Chen et al., 2024b], and has since been applied to various LLM fields [Yuan et al., 2024a, Wu et al., 2024b, Gao et al., 2024a]. The core idea of self-play fine-tuning is to progressively refine the LLM by competing against itself. Formally, it can be framed as a two-player game between a main player and an opponent player, each of which corresponds to an instance policy of the same LLM. In this game, the main player aims to distinguish expert-annotated responses from generated ones, while the opponent player strives to produce high-quality synthetic responses that can deceive the main player.

Based on this framework, Chen et al. [2024b] introduce the seminal method named SPIN, which is designed to optimize the relative reward advantages of annotated responses over generated ones, and employs the policy from previous iteration as a reference for reward computations. However, with the evolution of LLM, the relative advantage between two types of data gradually diminishes over time. When the relative advantage is eliminated, the optimization objective of SPIN degenerates into a constant independent of the policy, allowing *any policy* to become an optimal solution and thus leading to *unstable optimization*. Moreover, the incorporation of reference policy in SPIN causes a *misalignment* between the reward formulation used for training and the log likelihood that guides response generation. Empirical evidences in Figure 3 reveal that although annotated responses yield higher rewards than generated responses, their log likelihoods become lower than synthetic samples during iterations. In other words, higher rewards do not necessarily lead to higher generation priority.

To address these limitations, we propose a new self-play fine-tuning method, termed Triplet-based Self-Play fIne-tuNing (T-SPIN). The basic intuition is to incorporate historical advantages for stable evolution, and eliminate the reference policy for training-generation alignments. Specifically, we introduce two key modifications to the self-play framework. First, the main player is required not only to distinguish between real and synthetic data, but also discriminate *proto-synthetic* samples produced by the initial policy. By incorporating triplet inputs of real, synthetic and proto-synthetic samples, T-SPIN optimizes current and historical advantages simultaneously. Even if the current advantage diminishes (i.e., synthetic samples become indistinguishable from annotated ones), the historical advantages between synthetic responses and proto-synthetic responses still exist, ensuring a stable optimization over the policy model. Second, the opponent player maximizes the confidence of synthetic responses with an *entropy constraint*, where the confidence measures how likely the main player believes a response comes from the annotated data. Through theoretical analysis, we demonstrate that the incorporation of the entropy constraint removes the need for reference policy, aligning the principles of training and generation. Empirical results show that T-SPIN consistently maintains higher log-likelihoods for annotated data compared to synthetic samples.

We conduct experiments on Zephyr-7B [Tunstall et al., 2023] and Mistral-7B [Jiang et al., 2023], and select 50k annotated samples from Ultrachat200k [Ding et al., 2023] as the training set. The performance is evaluated across diverse tasks to assess multi-dimensional capabilities of LLMs, including math and logic, multi-domain knowledge, commonsense reasoning, and instruction-following. Experimental results demonstrate that T-SPIN significantly improves performance, particularly in math and instruction-following, and achieves superior and stable performances during iterations compared to SPIN. Furthermore, we also show that T-SPIN, using only 25% annotated samples, can achieve comparable performance (or even better on certain tasks) to supervised fine-tuning with the full training set, highlighting the effectiveness of T-SPIN in addressing the scarcity of annotated data.

## 2 Preliminary

In this section, we provide a brief overview of three fine-tuning strategies (i.e., supervised fine-tuning, self-play fine-tuning, and our T-SPIN) for LLMs. The comparisons are also illustrated in Figure 1.

#### 2.1 Supervised fine-tuning

A common strategy to adapt a pretrained model to downstream tasks is fine-tuning it on high-quality annotated data. Specifically, given a prompt  $\mathbf{x}$  and its corresponding annotated response  $\mathbf{y}$ , supervised fine-tuning (SFT) aims to align the model with  $(\mathbf{x}, \mathbf{y})$ . Formally, the goal of SFT is to learn a policy  $\pi_{\text{SFT}}$  by maximizing the log-likelihood of high-quality annotated data:

$$\pi_{SFT} = \operatorname*{argmax}_{\pi_{\theta}} \mathbb{E}_{\mathbf{x} \sim q(\cdot), \mathbf{y} \sim \pi_{data}(\cdot | \mathbf{x})} \left[ \log \pi_{\theta}(\mathbf{y} | \mathbf{x}) \right], \tag{1}$$

where  $q(\cdot)$  and  $\pi_{data}(\cdot|\mathbf{x})$  denote the distributions of prompts and annotated responses, respectively. Mathematically, (1) is equivalent to the solution of minimizing the KL divergence between the target distribution  $\pi_{data}(\cdot|\mathbf{x})$  and the policy  $\pi_{\theta}(\cdot|\mathbf{x})$ :

$$\pi_{\text{SFT}} = \operatorname*{argmin}_{\pi_{\theta}} \mathbb{E}_{\mathbf{x} \sim q(\cdot)} \left[ \mathbb{D}_{\text{KL}} (\pi_{data}(\cdot | \mathbf{x}) || \pi_{\theta}(\cdot | \mathbf{x})) \right]. \tag{2}$$

It can be verified that the solution (2) exactly matches the target distribution of annotated responses, i.e.,  $\pi_{SFT}(\cdot|\mathbf{x}) = \pi_{data}(\cdot|\mathbf{x})$ . Unfortunately, achieving the optimal policy typically demands extensive high-quality samples, resulting in high costs for data curation and annotation, and thereby limiting the practical applicability of supervised fine-tuning [Wang et al., 2024a, Tan et al., 2024].

### 2.2 Self-play fine-tuning

Self-play fine-tuning has emerged as a promising way to address the scarcity of expert-annotated data when adapting LLMs to downstream tasks [Chen et al., 2024b, Yuan et al., 2024a]. Different from SFT that runs over a fixed annotated dataset, self-play fine-tuning employs both annotated samples and self-generated responses as inputs, with the latter continuously refined as the model evolves.

Specifically, at each iteration t+1, the model first generates a synthetic response  $\mathbf{y}'$  for a given prompt  $\mathbf{x}$  by sampling from the policy  $\pi_{\theta_t}(\cdot|\mathbf{x})$ . Then, the model utilizes the prompt  $\mathbf{x}$  and pairwised samples (i.e., the annotated response  $\mathbf{y}$  and the synthetic response  $\mathbf{y}'$ ) to obtain the policy  $\pi_{\theta_{t+1}}$ . Let the reward for a response be formulated by the log ratio of the likelihood between the current policy  $\pi_{\theta}$  and the previous policy  $\pi_{\theta_t}$  [Rafailov et al., 2023]. Then, the objective of SPIN [Chen et al., 2024b] is to maximize the gap (referred to as the relative advantage) between the reward of annotated response  $\mathbf{y}$  and that of synthetic response  $\mathbf{y}'$ :

$$\pi_{\theta_{t+1}} = \operatorname*{argmin}_{\pi_{\theta}} \mathbb{E} \left[ \ell \left( \lambda \log \frac{\pi_{\theta}(\mathbf{y}|\mathbf{x})}{\pi_{\theta_{t}}(\mathbf{y}|\mathbf{x})} - \lambda \log \frac{\pi_{\theta}(\mathbf{y}'|\mathbf{x})}{\pi_{\theta_{t}}(\mathbf{y}'|\mathbf{x})} \right) \right], \tag{3}$$

where  $\lambda>0$  denotes a hyper-parameter,  $\ell(\cdot)$  is a convex and monotonically decreasing function, and the expectation operation is taken over the distributions  $\mathbf{x}\sim q(\cdot)$ ,  $\mathbf{y}\sim \pi_{data}(\cdot|\mathbf{x})$  and  $\mathbf{y}'\sim \pi_{\theta_t}(\cdot|\mathbf{x})$ . In the following, we elaborate two major limitations of SPIN: *unstable optimization* during iterations and *misalignment* between training and generation.

**Unstable optimization.** The instability of SPIN arises from the inherent structure of the objective in (3). To be precise, we consider the case where the synthetic response  $\mathbf{y}'$  becomes identical to the annotated response  $\mathbf{y}$  (i.e.,  $\mathbf{y}' = \mathbf{y}$ ). In this case, the advantage of  $\mathbf{y}$  over  $\mathbf{y}'$  vanishes, causing the objective (3) to degenerate into a constant that is independent of the policy  $\pi_{\theta}$ . As a result, *any* policy  $\pi_{\theta}$  becomes optimal to (3), leading to unstable optimization or even performance degradation.

Misalignment between training and generation. Given the reference policy  $\pi_{\theta_t}$  from the iteration t, the objective (3) is designed to optimize the reward gap between two types of responses, where the reward for a response  $\hat{\mathbf{y}}$  is defined as  $r(\mathbf{x}, \hat{\mathbf{y}}) = \lambda \log \pi_{\theta}(\hat{\mathbf{y}}|\mathbf{x}) - \lambda \log \pi_{\theta_t}(\hat{\mathbf{y}}|\mathbf{x})$ . However, it is important to note that the formulation of  $r(\mathbf{x}, \hat{\mathbf{y}})$  differs from the principle used for generation, i.e.,  $\log \pi_{\theta}(\hat{\mathbf{y}}|\mathbf{x})$ . In other words, due to the incorporation of reference policy  $\pi_{\theta_t}$ , a higher reward, i.e.,  $r(\mathbf{x}, \mathbf{y}) > r(\mathbf{x}, \mathbf{z})$  for two responses  $\mathbf{y}$  and  $\mathbf{z}$ , does not necessarily imply a higher generation probability, i.e.,  $\log \pi_{\theta}(\mathbf{y}|\mathbf{x}) > \log \pi_{\theta}(\mathbf{z}|\mathbf{x})$  may not hold. This demonstration has been empirically verified by our experiments, as illustrated in Figure 3.

### 3 Triplet-based self-play fine-tuning

In this section, we introduce our triplet-based self-play fine-tuning (T-SPIN), a stable and reference-free method for fine-tuning LLMs. Building upon the self-play framework, our T-SPIN involves two players: a main player and an opponent player, both of which are different policy instances from the same LLM. For clarity, we describe the update rules for two players separately.

The update rule for main player. Recall that the goal of the main player is to distinguish among the annotated response  $\mathbf{y}$ , the synthetic response  $\mathbf{y}'$  generated from the previous policy  $\pi_{\theta_t}$  and the proto-synthetic response  $\mathbf{y}_0$  produced by the initial policy  $\pi_{\theta_0}$ . Conceptually, the triple responses are expected to satisfy the relation  $\mathbf{y} \succeq \mathbf{y}' \succeq \mathbf{y}_0$ , where the operator  $\succeq$  indicates that the main player has greater preference for one response than in another. In other words,  $\mathbf{y}_0$  serves as an anchor that helps to evaluate how far the current policy  $\pi_{\theta_t}$  has progressed. This design is partially inspired by the schema-equilibration process in cognitive development, where knowledge builds upon and refines pre-existing cognitive anchors [Piaget, 1952]. Formally, we are motivated by the integral probability metric (IPM) [Müller, 1997] to establish the objective of the main player. In detail, let  $c(\mathbf{x}, \hat{\mathbf{y}})$  denote the confidence with which the main player believes a response  $\hat{\mathbf{y}}$  is sampled from the real-world data distribution  $\pi_{\text{data}}(\cdot|\mathbf{x})$ . Then, the goal of main player is to find a confidence function

$$c_{t+1} = \underset{c \in \mathcal{C}}{\operatorname{argmin}} \mathbb{E} \left[ \ell \left( c(\mathbf{x}, \mathbf{y}) - c(\mathbf{x}, \mathbf{y}') \right) + \beta \ell \left( c(\mathbf{x}, \mathbf{y}') - c(\mathbf{x}, \mathbf{y}_0) \right) \right], \tag{4}$$

where  $\beta \geq 0$  denotes a trade-off hyper-parameter, and  $\mathcal{C}$  denotes a function class that will be specified later, and  $\ell(\cdot)$  denotes a monotonically decreasing function, e.g., the logistic function  $\ell(x) = -\log \sigma(x)$  and the exponential function  $\ell(x) = \exp(-x)$ . The expectation operation here is taken over the distributions  $\mathbf{x} \sim q(\cdot)$ ,  $\mathbf{y} \sim \pi_{data}(\cdot|\mathbf{x})$ ,  $\mathbf{y}' \sim \pi_{\theta_t}(\cdot|\mathbf{x})$  and  $\mathbf{y}_0 \sim \pi_{\theta_0}(\cdot|\mathbf{x})$ . Given the triple responses  $\mathbf{y}$ ,  $\mathbf{y}'$ , and  $\mathbf{y}_0$ , the objective in (4) consists of two terms:

- The first term  $\ell\left(c(\mathbf{x},\mathbf{y})-c(\mathbf{x},\mathbf{y}')\right)$  measures the gap between the real response  $\mathbf{y} \sim \pi_{\text{data}}(\cdot|\mathbf{x})$  and the synthetic one  $\mathbf{y}' \sim \pi_{\theta_t}(\cdot|\mathbf{x})$ , reflecting the relative advantages of  $\pi_{\text{data}}$  over  $\pi_{\theta_t}$ ;
- The second term  $\ell\left(c(\mathbf{x},\mathbf{y}')-c(\mathbf{x},\mathbf{y}_0)\right)$  compares the synthetic response  $\mathbf{y}'$  to the proto-synthetic response  $\mathbf{y}_0 \sim \pi_{\theta_0}(\cdot|\mathbf{x})$ , capturing historical advantages of  $\pi_{\theta_t}$  over  $\pi_{\theta_0}$  during iterations.

During the self-play process, the opponent player evolves over time, generating synthetic responses  $\mathbf{y}'$  that progressively approximate the real response  $\mathbf{y}$ . As a result, the current advantage, i.e., the first term in (4), gradually diminishes. In contrast, since  $\mathbf{y}_0$  remains fixed during iterations, the historical advantage, i.e., the second term in (4), remains meaningful even when the current advantage vanishes, ensuring that (4) does not degenerate to a constant independent of  $\pi_{\theta}$ .

The update rule for opponent player. In this part, we suppose that the main player has learned a confidence function  $c_{t+1}$  capable of effectively distinguishing between  $\mathbf{y}$  and  $\mathbf{y}'$ . Then, we proceed to elaborate on the update rule for the opponent player. Note that the goal of the opponent player is to generate a response  $\mathbf{y}'$  that misleads the main player to the best of its ability. Therefore, with the given confidence function  $c_{t+1}$ , the opponent player aims to find a policy  $\pi_{\hat{\theta}}$  that maximizes the expected confidence score over synthetic responses  $\mathbf{y}'$  generated from itself:

$$\pi_{\hat{\theta}} = \operatorname*{argmax}_{\pi_{\theta}} \mathbb{E}_{\mathbf{x} \sim q(\cdot), \mathbf{y}' \sim \pi_{\theta}(\cdot | \mathbf{x})} \left[ c_{t+1}(\mathbf{x}, \mathbf{y}') \right] + \alpha \mathbb{E}_{\mathbf{x} \sim q(\cdot)} \left[ \mathcal{H}(\pi_{\theta}(\cdot | \mathbf{x})) \right], \tag{5}$$

where  $\alpha > 0$  denotes the regularization parameter, and  $\mathcal{H}(\cdot)$  denotes the entropy constraint. In (5), the first part reflects the confidence that how the main player believes that  $\mathbf{y}'$  is sampled from the real-world distribution  $\pi_{data}$ , and the second part serves to stabilize the optimization. It is worth noting that (5) has a closed-form solution, as stated below.

**Proposition 1.** Given a general confidence function  $c_{t+1}$ , the optimal policy of the opponent player in (5) takes the form of

$$\pi^*(\mathbf{y}'|\mathbf{x}) = \frac{\exp(c_{t+1}(\mathbf{x}, \mathbf{y}')/\alpha)}{\sum_{\mathbf{z} \in \mathcal{Y}} \exp(c_{t+1}(\mathbf{x}, \mathbf{z})/\alpha)}$$

where  $\mathcal{Y}$  denotes the set that contains all possible responses for the prompt  $\mathbf{x}$ .

Note that due to the generality of  $c_{t+1}$ , the closed-form solution  $\pi^*(\mathbf{y}'|\mathbf{x})$  may not be attainable in the policy space of LLM. To resolve this issue, we need to choose a specific form of the function class  $\mathcal{C}$  such that the closed-form solution becomes *tractable*. Specifically, we hope to find a policy

### Algorithm 1 Triplet-based Self-Play Fine-Tuning (T-SPIN)

**Inputs**: A set of high-quality data  $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n$ , a pretrained LLM  $\pi_{\theta_0}$ , and the number of iterations T

- 1: Produce the proto-synthetic response  $y_0$  with the policy  $\pi_{\theta_0}$  for each prompt x.
- 2: **for**  $t = 0, 1, \dots, T 1$  **do**
- 3: For each prompt x, generate one synthetic response y' by the policy  $\pi_{\theta_t}$ .
- 4: Update the policy  $\pi_{\theta_{t+1}}$  by minimizing (7) with prompt x and triplet responses  $\{y, y', y_0\}$ .
- 5: end for

 $\pi_{\theta}$  that satisfies  $\pi_{\theta}(\mathbf{y}'|\mathbf{x}) = \pi^*(\mathbf{y}'|\mathbf{x})$  with the confidence  $c_{t+1}$  learned from a certain class C. To this end, according to Proposition 1, we can obtain  $\pi_{\theta}(\mathbf{y}'|\mathbf{x}) \propto \exp(c_{t+1}(\mathbf{x},\mathbf{y}')/\alpha)$ , which implies  $c_{t+1}(\mathbf{x},\mathbf{y}') = \alpha \log \pi_{\theta}(\mathbf{y}'|\mathbf{x})$ . In other words, we can choose the function class

$$C = \{\alpha \log \pi_{\theta}(\cdot | \mathbf{x}) | \theta \in \Theta\}$$
 (6)

where  $\Theta$  denote the parameter space of LLM, and optimize (4) to obtain the confidence  $c_{t+1}(\mathbf{x}, \mathbf{y}') = \alpha \log \pi_{\theta}(\mathbf{y}'|\mathbf{x})$  where the policy  $\pi_{\theta}$  is exactly the optimal solution of (5). Therefore, we do not need to solve (5) explicitly. Moreover, since (6) depends solely on the policy  $\pi_{\theta}$  and does not rely on *any* reference policy, it inherently avoids the training-generation misalignment issue suffered by SPIN.

**The end-to-end objective.** We consolidate the two update rules into a unified end-to-end objective by substituting (6) into (4), and obtain the loss function for T-SPIN as follows:

$$\mathcal{L}_{\text{T-SPIN}}(\theta) = \mathbb{E}\left[\ell\left(\alpha \log \pi_{\theta}(\mathbf{y}|\mathbf{x}) - \alpha \log \pi_{\theta}(\mathbf{y}'|\mathbf{x})\right) + \beta \ell\left(\alpha \log \pi_{\theta}(\mathbf{y}'|\mathbf{x}) - \alpha \log \pi_{\theta}(\mathbf{y}_{0}|\mathbf{x})\right)\right]. \tag{7}$$

The expectation is taken over distributions  $\mathbf{x} \sim q(\cdot), \mathbf{y} \sim \pi_{data}(\cdot|\mathbf{x}), \mathbf{y}' \sim \pi_{\theta_t}(\cdot|\mathbf{x}), \mathbf{y}_0 \sim \pi_{\theta_0}(\cdot|\mathbf{x})$ . Overall, the training of T-SPIN can be viewed as alternating updates between the main player and the opponent player. At each iteration t+1, the opponent player first generates synthetic responses  $\mathbf{y}'$  sampled from the previous iteration policy  $\pi_{\theta_t}$ . Then, the main player takes the triplet responses  $\mathbf{y}, \mathbf{y}'$  and  $\mathbf{y}_0$  as inputs to obtain its confidence  $c_{t+1} = \alpha \log \pi_{\theta_{t+1}}$ . After that, the opponent player chooses the policy  $\pi_{\theta_{t+1}}$  for the next iteration. The detailed steps are presented in Algorithm 1.

Comparisons between T-SPIN and SPIN. Both methods are designed to handle the limited available data when adapting LLMs to downstream tasks, but they fundamentally differ in (i) optimization principles, (ii) implicit reward formulation, and (iii) computational overhead, as detailed below:

- SPIN only optimizes the current advantage between the annotated response  $\mathbf{y} \sim \pi_{data}(\cdot|\mathbf{x})$  and the synthetic response  $\mathbf{y}' \sim \pi_{\theta_t}(\cdot|\mathbf{x})$  in (3), which may result in the instability issue during iterations. In contrast, T-SPIN optimizes not only the current advantage of  $\mathbf{y}$  over  $\mathbf{y}'$ , but also the historical advantage of  $\mathbf{y}'$  over the proto-synthetic response  $\mathbf{y}_0 \sim \pi_{\theta_0}(\cdot|\mathbf{x})$ , ensuring a stable evolution and thereby mitigating the unstable optimization;
- In SPIN, the reward function for a response  $\mathbf{z}$  is formulated as  $r_{\theta}(\mathbf{z}|\mathbf{x}) = \lambda \log \pi_{\theta}(\mathbf{z}|\mathbf{x}) \lambda \log \pi_{\theta_t}(\mathbf{z}|\mathbf{x})$ . This formulation misaligns with the generation metric (i.e., the log-likelihood), resulting in potential suboptimal performance. By contrast, the reward function in T-SPIN (also referred to as the confidence function under the self-play framework) is defined as  $r_{\theta}(\mathbf{z}|\mathbf{x}) = \alpha \log \pi_{\theta}(\mathbf{z}|\mathbf{x})$ , which matches the generation metric and thus resolves the misalignment issue;
- Compared to SPIN, T-SPIN incurs only a marginal increase in computational and memory costs. Specifically, T-SPIN performs a one-time generation of  $\mathbf{y}_0$  at the beginning, which is not repeated in subsequent iterations. In practice, the costs of this generation is negligible compared to the iterative training process. Moreover, while T-SPIN expands each training instance from a pair of  $(\mathbf{y}, \mathbf{y}')$  to a triplet of  $(\mathbf{y}, \mathbf{y}', \mathbf{y}_0)$ , the overall memory costs remains nearly unchanged, as it eliminates the need for a reference model, which is necessary in SPIN. Consequently, the additional overhead of T-SPIN is minimal compared to SPIN.

**Gradient analysis.** To mechanistically understand the impact of T-SPIN on the training process, we examine the gradient of (7) with respect to the parameter  $\theta$ , as presented in the following theorem.

**Theorem 1.** Define  $u(\pi_{\theta}) = \log \pi_{\theta}(\mathbf{y}|\mathbf{x}) - \log \pi_{\theta}(\mathbf{y}'|\mathbf{x})$  and  $v(\pi_{\theta}) = \log \pi_{\theta}(\mathbf{y}'|\mathbf{x}) - \log \pi_{\theta}(\mathbf{y}_{0}|\mathbf{x})$ . Then, given a monotonically decreasing function  $\ell(x)$ , the gradient of (7) with respect to  $\theta$  is

$$\nabla_{\theta} \mathcal{L}_{\text{T-SPIN}}(\theta) = \alpha \mathbb{E} \Big[ \ell'(\alpha u(\pi_{\theta})) \cdot (\nabla_{\theta} \log \pi_{\theta}(\mathbf{y}|\mathbf{x}) - \nabla_{\theta} \log \pi_{\theta}(\mathbf{y}'|\mathbf{x})) \\ + \beta \ell'(\alpha v(\pi_{\theta})) \cdot (\nabla_{\theta} \log \pi_{\theta}(\mathbf{y}'|\mathbf{x}) - \nabla_{\theta} \log \pi_{\theta}(\mathbf{y}_{0}|\mathbf{x})) \Big]$$
(8)

Table 1: Performance (%) comparisions among three fine-tuning methods, where SPIN and T-SPIN are trained on 50k annotated data from Ultrachat200k, and SFT is trained on the complete dataset. At the average (Avg) column, we also report improvements (or degradations) over the previous iteration. Abbreviations: MMLU-Pro  $\rightarrow$  MMLUP, Winograde  $\rightarrow$  WG, Instruction Following  $\rightarrow$  IF.

Model		Math & Logic			Multi-Domain Knowledge   Common				ense Ro	easoning	IF	Avg
		GSM8K	MATH	MUSR	MMLU	MMLUP	GPQA	HellaSwag	WG	BBH	IFEval	
Zep	hyr-7B	25.85	1.75	38.92	56.90	28.88	28.91	82.79	74.19	44.60	2.76	38.56
5	SFT	42.25	3.10	39.18	57.29	29.19	28.28	83.44	73.16	44.90	19.31	42.01
	Iter0	30.06	4.46	41.67	56.25	28.39	29.53	84.29	73.56	44.41	8.12	40.07
z	Iter1	32.47	3.05	39.02	56.86	27.88	29.06	83.76	73.72	44.59	7.95	$39.84_{(-0.23)}$
SPIN	Iter2	36.62	2.57	41.27	57.74	28.67	28.39	83.75	73.95	43.95	14.46	$41.14_{(+1.30)}$
01	Iter3	33.32	2.75	39.41	55.88	27.98	27.61	82.64	74.19	43.19	22.25	$40.92_{(-0.22)}$
	Iter4	35.54	2.72	38.21	53.59	25.95	26.21	83.48	74.51	43.10	22.88	$40.62_{(-0.30)}$
-	Iter0	36.20	3.29	37.05	56.42	28.88	28.05	82.97	73.48	44.17	7.02	39.75
IN	Iter1	40.64	3.68	37.84	56.89	29.13	28.92	82.95	73.72	44.82	26.98	$42.56_{(+2.81)}$
	Iter2	40.69	3.63	38.91	56.89	29.26	29.40	83.37	73.64	44.48	27.67	$42.79_{(+0.23)}$
Ė	Iter3	40.71	3.82	39.18	57.67	29.37	30.12	83.30	73.74	45.27	29.15	$43.23_{(+0.44)}$
	Iter4	40.67	3.84	39.72	57.68	29.32	30.44	83.12	73.80	45.00	31.08	$43.47_{(+0.24)}$
45.0		T-SPIN			40.0	T-SPIN	.1		45.0	T-SF	OTNI	
41	1.0	SPIN			30.0	SPIN	N		43.0	SPI		

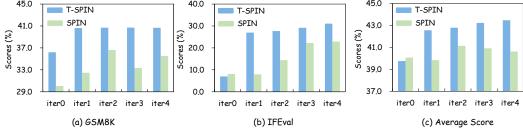


Figure 2: Performance (%) comparisons between T-SPIN and SPIN on two tasks: GMS8K and IFEval over 5 iterations. The average scores over 10 different tasks are also illustrated in the right panel.

where  $\ell'(\cdot)$  denotes the derivative of  $\ell(\cdot)$ , and the expectation is taken over the distributions  $\mathbf{x} \sim q(\cdot), \mathbf{y} \sim \pi_{data}(\cdot|\mathbf{x}), \mathbf{y}' \sim \pi_{\theta_t}(\cdot|\mathbf{x}), \mathbf{y}_0 \sim \pi_{\theta_0}(\cdot|\mathbf{x})$ . With a specific choice of  $\ell(x)$ , we can derive more explicit forms of the gradient. Detailed derivation are provided in Appendix A.2.

Note that the function  $\ell(x)$  is monotonically decreasing, leading to a non-positive derivative  $\ell'(x) \leq 0, \forall x \in \mathbb{R}$ . Therefore, during the training, the gradient in (8) tends to increase the likelihood of  $\mathbf{y}$  and decrease the likelihood of  $\mathbf{y}_0$  with weights corresponding to the current advantage  $u(\pi_\theta)$  and the historical advantage  $v(\pi_\theta)$ , respectively. For the synthetic response  $\mathbf{y}'$ , the update direction is influenced by both current and historical advantages, with the latter weighted by hyperparameter  $\beta$ .

### 4 Experiments

In this section, we commence by conducting empirical evaluations of T-SPIN, SPIN, and SFT on a set of benchmarks with 10 different tasks. Then, we provide in-depth ablation analysis to investigate the impact of components (i.e., historical advantages, reward function and hyperparameters) in T-SPIN.

### 4.1 Experimental setups

Models and training settings. Following the setup in Chen et al. [2024b], we randomly sample 50k prompt-response pairs from the first round of interactions in Ultrachat200k [Ding et al., 2023] as the training set, and choose Zephyr-7B-SFT-Full [Tunstall et al., 2023] and Mistral-7B-v0.1 [Jiang et al., 2023] as pretrained LLMs in our experiments. We employ the commonly used template: "### Instruction: {prompt}\n\n### Response:" [Taori et al., 2023] to construct the prompts, and set  $\ell(x) = -\log \sigma(x)$  by the logistic loss function in the objectives of T-SPIN and SPIN. The

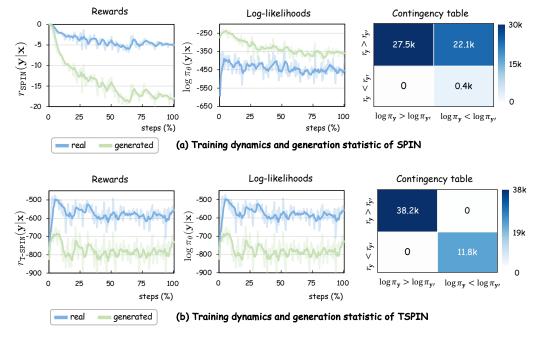


Figure 3: Comparisons between SPIN and T-SPIN at iteration 1: (a) Training dynamics (including rewards and log-likelihoods of  $\mathbf{y}$  and  $\mathbf{y}'$ ) and generation statistic (i.e., associations between rewards and log-likelihoods) measured on the training set for iteration 1 of SPIN; (b) Training dynamics and generation statistic of T-SPIN. For brevity, we denote  $r_{\mathbf{y}} = r(\mathbf{y}|\mathbf{x})$  and  $\log \pi_{\mathbf{y}} = \log \pi(\mathbf{y}|\mathbf{x})$ .

training process adheres to the iterative framework described in Section 3. At each iteration, we first employ the latest policy to generate synthetic responses, which are then combined with annotated responses as the training data for the next iteration. In total, there are 50k synthetic responses and 50k annotated responses in each iteration. For T-SPIN, we additionally incorporate 50k proto-synthetic responses produced by the initial policy, into the training set.

**Evaluation benchmarks.** We evaluate the performance across diverse tasks, each designed to assess specific capabilities of LLMs. These tasks include mathematical and logical reasoning with GSM8K [Cobbe et al., 2021], MATH [Hendrycks et al., 2021b], MUSR [Sprague et al., 2024]; multi-domain knowledge understanding through MMLU [Hendrycks et al., 2021a], MMLU-Pro [Wang et al., 2024d], GPQA [Rein et al., 2023]; commonsense reasoning by HellaSwag [Zellers et al., 2019], Winogrande [Sakaguchi et al., 2021] and BBH [Suzgun et al., 2022]; and instruction-following via IFEval [Zhou et al., 2023]. For each task, we follow the standard evaluation protocols in the literature. For instance, we utilize 3-shot examples and normalized accuracy as the metric in BBH. More detailed introductions about the evaluation benchmarks can be found in Appendix B.2.

#### 4.2 Experimental results

Comparisions to SPIN and SFT. We first compare T-SPIN with SPIN and SFT on multiple tasks, where two self-play methods are trained with only 50k annotated data, while SFT is trained with 200k annotated data. Table 1 summarizes the results on Zephyr-7B, and the results on Mistral-7B are deferred to the Appendix B.3. The results demonstrate that T-SPIN significantly improve the performances of pretrained models, and outperform other two fine-tuning methods. Specifically, T-SPIN achieves a significant improvement in average score compared to Zephyr-7B, particularly excelling in GSM8K and IFEval tasks with 14.82-point and 28.32-point increases respectively, which is also illustrated in Figure 2(a) and Figure 2(b). Moreover, from the results shown in Figure 2(c), we also observe that SPIN suffers from performance instability over iterations, reaching the peak performance at the third iteration and then suffering the degradation at the subsequent iterations. In contrast, T-SPIN demonstrates stable evolutions over all iterations. Compared to SFT that utilizes the complete annotated set as training data, T-SPIN with only 50k annotated data achieves comparable (or better) performances, highlighting the effectiveness of T-SPIN when facing the data scarcity issue.

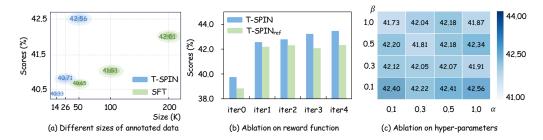


Figure 4: Performance comparisons over different settings: (a) comparison between T-SPIN and SFT with varying amounts of annotated data; (b) performances of T-SPIN versus the variant without historical advantage (w/o H-A); (c) robustness analysis with respect to hyperparameters  $\alpha$  and  $\beta$ .

**Aligning training rewards and generation log-likelihoods.** We then show that T-SPIN is able to alleviate the misalignment issue suffered by SPIN. We plot the training dynamics and generation statistics of SPIN and T-SPIN on Zephyr-7B at the iteration 1 in Figure 3(a) and Figure 3(b), respectively. We observe that in SPIN, annotated responses yield higher rewards than synthetic responses, but their log-likelihoods are lower than those of synthetic responses, shown in the left and middle panels of Figure 3(a). The reason lies in that the reward function of SPIN incorporates both the current policy  $\pi_{\theta}$  and the reference policy  $\pi_{\theta_t}$  from the previous iteration, expressed as  $r_{\text{SPIN}}(\mathbf{x}, \mathbf{z}) = \lambda \log \pi_{\theta}(\mathbf{z}|\mathbf{x}) - \lambda \log \pi_{\theta_{\star}}(\mathbf{z}|\mathbf{x})$  for a response  $\mathbf{z}$ . This formulation differs from the generation principle that solely relies on  $\log \pi_{\theta}(\mathbf{z}|\mathbf{x})$ . Consequently, responses with higher reward does not necessarily enjoy higher generation probability. To further illustrate this phenomenon, we also analyze the generation statistics of rewards and log-likelihoods at the right panel of Figure 3(a). The results indicate that nearly half of the samples exhibit higher rewards but lower log-likelihoods. By contrast, T-SPIN effectively mitigates the mismatch issue by aligning the reward formulation for training and log-likelihood used for generation, i.e.,  $r_{\texttt{T-SPIN}}(\mathbf{x}, \mathbf{z}) = \alpha \log \pi_{\theta}(\mathbf{z}|\mathbf{x})$  where  $\alpha = 1$  in this experiment. Empirical evidence, shown in the left and middle panels of Figure 3(b), demonstrates that both the rewards and log-likelihoods of real responses are consistently higher than those of generated responses. We also examine the generation statistic of two types of responses in terms of rewards and log-likelihoods, with similar trends observed at the right panel of Figure 3(b).

Fine-tuning with less annotated data. We now demonstrate that T-SPIN can achieve comparable or even superior performances to SFT, while using less annotated training data. To validate this claim, we conduct experiments with different sizes of annotated data. Specifically, we construct subsets of 14k, 26k, 50k, 100k, and 200k samples from Ultrachat200k [Ding et al., 2023], where each larger subset encompasses all smaller ones. For comparison, we train T-SPIN on Zephyr-7B using the smaller subsets (14k, 26k, and 50k samples) and SFT using the larger subsets (50k, 100k) and 200k samples). We report the performances of T-SPIN after 2 iterations and SFT after 2 consecutive epochs. Figure 4(a) illustrates the empirical results, from which we can observe that: (i) both fine-tuning methods show improved performance with the increased annotated data size; (ii) with the same amount of annotated data (i.e., 50k), T-SPIN achieves a superior average score of 42.56% compared to 40.65% of SFT; (iii) T-SPIN using only 50k annotated responses attains 42.56% average score, slightly outperforming 42.01% of SFT trained with 200k annotated data.

#### 4.3 Ablation studies

In this part, we conduct ablation studies on Zephyr-7B to investigate the impact of each component (i.e., historical advantages, reference-free reward function, and hyper-parameters) in T-SPIN.

**Historical advantages.** In T-SPIN, the historical advantage, i.e., the second term in (7), is introduced to address the instability issues caused by the vanishing current advantage, i.e., the first term in (7). To verify the effectiveness of historical advantages, we conduct ablation experiments by comparing T-SPIN with a variant that removes the historical advantage from (7), referred to as w/o H-A. Specifically, for T-SPIN, we initialize  $\beta=0$  in the first iteration, as the historical advantage is identical to the current advantage at iteration 0, and set  $\beta=0.1$  for all subsequent iterations. For w/o H-A, we fix  $\beta=0$  over all iterations. We report the performances of T-SPIN and w/o H-A on all tasks in Table 2. The average performances over different tasks reveal two key findings: (i) both T-SPIN and w/o H-A are able to improve the performances, although w/o H-A suffers from a performance decline

Table 2: Performance (%) comparisions between T-SPIN and w/o H-A on Zephyr-7B. At the average (Avg) column, we also report the performance improvements (or degradations) over the previous iteration. Abbreviations: MMLU-Pro  $\rightarrow$  MMLUP, Winograde  $\rightarrow$  WG, Instruction Following  $\rightarrow$  IF.

Methods	Math & Logic			Multi-E	Avg						
	GSM8K	MATH	MUSR	MMLU	MMLUP	GPQA	HellaSwag	WG	BBH	IFEval	
Iter0	36.20	3.29	37.05	56.42	28.88	28.05	82.97	73.48	44.17	7.02	39.75
Iter1	29.89	3.33	37.06	57.51	28.98	27.35	82.13	73.16	44.63	10.46	$39.45_{(-0.30)}$
₩ Iter2	38.95	3.19	37.83	56.72	27.65	28.19	83.29	73.95	44.64	21.87	$41.63_{(+2.18)}$
Ş Iter3	41.71	2.63	37.59	56.87	27.82	28.96	83.28	73.32	44.63	20.10	$41.64_{(+0.01)}$
Iter4	38.12	3.83	37.05	56.96	28.40	28.61	83.71	74.11	43.86	23.20	$41.79_{(+0.05)}$
_ Iter1	40.64	3.68	37.84	56.89	29.13	28.92	82.95	73.72	44.82	26.98	$42.56_{(+2.81)}$
Iter2	40.69	3.63	38.91	56.89	29.26	29.40	83.37	73.64	44.48	27.67	$42.79_{(+0.23)}$
Iter2	40.71	3.82	39.18	57.67	29.37	30.12	83.30	73.74	45.27	29.15	$43.23_{(+0.44)}$
Iter4	40.67	3.84	39.72	57.68	29.32	30.44	83.12	73.80	45.00	31.08	$43.47_{(+0.24)}$

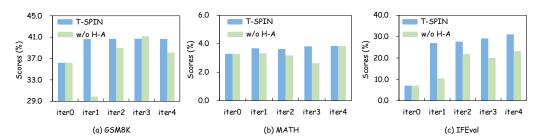


Figure 5: Performance (%) comparisons between T-SPIN and w/o H-A on GMS8K, MATH and IFEval.

at iteration 1; (ii) without the historical advantage, w/o H-A improves slowly and plateaus at 41.64% at iteration 3, whereas T-SPIN reaches a higher performance of 42.56% earlier, at iteration 1. When we dig into the performance of each task, we find that w/o H-A suffers from unstable performance trends, as shown in Figure 5. For example, on the task of GMS8K, w/o H-A suffers from a significant performance decline at the iteration 1. On MATH and IFEval, the performance degradation occurs at iteration 3. In contrast, T-SPIN demonstrates consistent improvements across all tasks and maintains stability when reaching the performance peak, highlighting the effectiveness of historical advantages.

**Reference-free reward function.** To address the training-generation misalignment issue in SPIN, T-SPIN introduce an entropy constraint into the self-play framework as shown in (5), from which we theoretically deliver a reference-free reward function for T-SPIN. To investigate its effectiveness, we compete our T-SPIN with a variant, namely T-SPIN<sub>ref</sub>, that still optimizes with triplet inputs, but uses the reward function with reference policy. The loss function of T-SPIN<sub>ref</sub> is given by:

$$\mathcal{L}_{\text{T-SPIN}_{\text{ref}}}(\theta) = \mathbb{E}\left[\ell\left(\alpha\log\frac{\pi_{\theta}(\mathbf{y}|\mathbf{x})}{\pi_{\theta_{t}}(\mathbf{y}|\mathbf{x})} - \alpha\log\frac{\pi_{\theta}(\mathbf{y}'|\mathbf{x})}{\pi_{\theta_{t}}(\mathbf{y}'|\mathbf{x})}\right) + \beta\ell\left(\alpha\log\frac{\pi_{\theta}(\mathbf{y}'|\mathbf{x})}{\pi_{\theta_{t}}(\mathbf{y}'|\mathbf{x})} - \alpha\log\frac{\pi_{\theta}(\mathbf{y}_{0}|\mathbf{x})}{\pi_{\theta_{t}}(\mathbf{y}_{0}|\mathbf{x})}\right)\right]. \tag{9}$$

We report the average performances of T-SPIN and T-SPIN<sub>ref</sub> on all tasks in Figure 4(b). From the results, we observe that T-SPIN demonstrates superior performance compared to T-SPIN<sub>ref</sub> over all iterations. This can be attributed to that using the reference-free reward function preserves the alignment between training and generation, which in turn facilitates performance improvement.

**Parameter robustness.** T-SPIN involves two hyper-parameters: (i) the trade-off parameter  $\beta$  that balances the current and historical advantages, and (ii) the regularization coefficient  $\alpha$  that controls the entropy penalty. To investigate the robustness of these parameters, we conduct experiments at iteration 1 with varying  $\alpha$  and  $\beta$ . We select values for  $\alpha$  and  $\beta$  from the set  $\{0.1, 0.3, 0.5, 1.0\}$ , and report the average performances over all 10 tasks with different parameters in Figure 4(c). The results demonstrate that T-SPIN is relatively insensitive to a wide range of  $\alpha$  and  $\beta$  with diverse values. Accordingly, we use  $\alpha = 1.0$  and  $\beta = 0.1$  as the default values in our experiments.

### 5 Related work

Self-play fine-tuning for LLMs. The self-play learning framework originates from TD-Gammon [Tesauro, 1995], and has produced profound impacts in various fields, such as multi-agent reinforcement learning [Lanctot et al., 2017, Muller et al., 2019, Bansal et al., 2018] and game-playing [Silver et al., 2017a,b, Anthony et al., 2017]. The central idea of this framework is to train the model by playing against its previous version for progressive improvements. Recently, self-play fine-tuning has been introduced into the investigations for LLMs [Chen et al., 2024b], which then leads to a series of subsequent works [Ren et al., 2024, Guo et al., 2024b, Yuan et al., 2024b, Rosset et al., 2024, Cheng et al., 2024, Pang et al., 2024, Wu et al., 2024b, Wang et al., 2025]. For example, Gao et al. [2024a] introduce the self-play mechanism into the recommendation systems to calibrate item-level predictions and suppress the bias issue, and Yuan et al. [2024a] incorporate the idea of self-play into diffusion models, resulting in a novel diffusion-based method. Different from the above studies that apply self-play fine-tuning to specific domains of LLMs, our work follows the seminal study of Chen et al. [2024b], investigating the fundamental nature of self-play fine-tuning, and aims to resolve the unstable evolution and misalignment issues in the original method.

**Data generation in LLMs.** Large language models are typically data-driven, heavily relying on the high-quality annotated data [Ouyang et al., 2022, Tan et al., 2024]. However, the available data is often limited for LLM training, and the costs of human expert annotation are prohibitively expensive. For this reason, recent efforts have been made to employ advanced LLMs, such as GPT series [OpenAI, 2023, 2024, 2025], Gemini series [Google, 2023, 2024, 2025], and Claude series [Anthropic, 2024a,b, 2025], to generate synthetic data for LLM fine-tuning [Josifoski et al., 2023, Gupta et al., 2023, Long et al., 2024, Qin et al., 2025]. It has been shown that in most domains, advanced expert LLMs are capable of generating high-quality data that can benefit others [Ubani et al., 2023, Peng et al., 2023, Xu et al., 2023, Yang et al., 2024, Choi et al., 2024, Zhang, 2024, Wang et al., 2024b, DeepSeek-AI et al., 2025]. Compared to these studies, our self-play fine-tuning method focuses on using data generated by the current policy model to train itself to achieve the progressive performance improvements, and offers an alternative to expert-LLM-based strategies.

### 6 Conclusion and discussion

In this paper, we investigate self-play fine-tuning for LLMs and identify two issues of the existing method SPIN: unstable optimization during iterations and misalignment between training and generation. To address these issues, we propose a novel self-play fine-tuning method named T-SPIN with two key designs. First, our T-SPIN simultaneously optimizes the current advantage of the annotated response y over the newly generated response y', and the historical advantage of y' over the proto-synthetic response  $y_0$ . When current advantages vanish, historical advantages can still guide the fine-tuning process, avoiding the unstable optimization. Second, T-SPIN incorporates the entropy constraint into the self-play framework, which eliminates the need for a reference policy in the reward computations, thus aligning the reward formulation and the generation metric. Empirically, we conduct extensive experiments and examine the performances of T-SPIN on a variety of tasks. Experimental results show that T-SPIN maintains stable evolution and achieves better performance than its contenders, and outperforms supervised fine-tuning by leveraging only 25% annotated data.

Limitations and future work. There are several promising directions for future research. First, at each iteration, T-SPIN needs to re-generate synthetic responses for all prompts, although some of them may no longer contribute to further evolutions. Consequently, a carefully designed filter strategy is necessary to choose valuable samples before generation. Second, in T-SPIN, the expert-annotated data is fixed during iterations. But, in real-world applications, the target distribution may shift over time, and hence the fixed annotated samples may not be suitable for online fine-tuning. For this reason, self-play fine-tuning for non-stationary target distributions is another interesting direction, for which we believe that some advanced techniques from online learning may offer valuable insights [Hazan, 2016, Orabona, 2019, Zhang et al., 2018, Wang et al., 2024c,e].

## Acknowledgments and Disclosure of Funding

This work was partially supported by NSFC (U23A20382), and the Collaborative Innovation Center of Novel Software Technology and Industrialization.

#### References

- J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, and et al. GPT-4 technical report. ArXiv e-prints, arXiv:2303.08774, 2023.
- T. Anthony, Z. Tian, and D. Barber. Thinking fast and slow with deep learning and tree search. In *Advances in Neural Information Processing Systems 30*, pages 5360–5370, 2017.
- Anthropic. Claude 3.5 sonnet. https://www.anthropic.com/news/claude-3-5-sonnet, 2024a.
- Anthropic. Introducing the next generation of claude. https://www.anthropic.com/news/claude-3-family, 2024b.
- Anthropic. Claude 3.7 sonnet and claude code. https://www.anthropic.com/news/claude-3-7-sonnet.2025.
- J. Bai, S. Bai, S. Yang, S. Wang, S. Tan, P. Wang, J. Lin, C. Zhou, and J. Zhou. Qwen-VL: A versatile vision-language model for understanding, localization, text reading, and beyond. *ArXiv e-prints*, arXiv:2308.12966, 2023.
- T. Bansal, J. Pachocki, S. Sidor, I. Sutskever, and I. Mordatch. Emergent complexity via multi-agent competition. In *Proceedings of the 6th International Conference on Learning Representations*, 2018.
- S. Chen, Y. Wang, Y.-F. Wu, Q.-G. Chen, Z. Xu, W. Luo, K. Zhang, and L. Zhang. Advancing tool-augmented large language models: Integrating insights from errors in inference trees. In *In Advances in Neural Information Processing Systems 37*, pages 106555–106581, 2024a.
- Z. Chen, Y. Deng, H. Yuan, K. Ji, and Q. Gu. Self-play fine-tuning convertsweak language models to strong language models. In *Proceedings of the 41st International Conference on Machine Learning*, pages 6621–6642, 2024b.
- P. Cheng, Y. Yang, J. Li, Y. Dai, T. Hu, P. Cao, N. Du, and X. Li. Adversarial preference optimization: Enhancing your alignment via RM-LLM game. In *Findings of the 62nd Conference of the Association for Computational Linguistics*, pages 3705–3716, 2024.
- J. Choi, J. Yun, K. Jin, and Y. Kim. Multi-News+: Cost-efficient dataset cleansing via llm-based data annotation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 15–29, 2024.
- K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano, C. Hesse, and J. Schulman. Training verifiers to solve math word problems. *ArXiv e-prints*, arXiv:2110.14168, 2021.
- H. Dai, Z. Liu, W. Liao, X. Huang, Y. Cao, Z. Wu, L. Zhao, S. Xu, W. Liu, N. Liu, S. Li, D. Zhu, H. Cai, L. Sun, Q. Li, D. Shen, T. Liu, and X. Li. AugGPT: Leveraging ChatGPT for text data augmentation. *ArXiv e-prints*, arXiv:2302.13007, 2023.
- DeepSeek-AI, D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi, X. Zhang, X. Yu, Y. Wu, Z. F. Wu, Z. Gou, Z. Shao, Z. Li, and Z. G. et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *ArXiv e-prints*, arXiv:2501.12948, 2025.
- N. Ding, Y. Chen, B. Xu, Y. Qin, Z. Zheng, S. Hu, Z. Liu, M. Sun, and B. Zhou. Enhancing chat language models by scaling high-quality instructional conversations. *ArXiv e-prints*, arXiv:2305.14233, 2023.
- C. Gao, R. Chen, S. Yuan, K. Huang, Y. Yu, and X. He. SPRec: Leveraging self-play to debias preference alignment for large language model-based recommendations. *ArXiv e-prints*, arXiv:2412.09243, 2024a.

- L. Gao, J. Tow, B. Abbasi, S. Biderman, S. Black, A. DiPofi, C. Foster, L. Golding, J. Hsu, A. L. Noac'h, H. Li, K. McDonell, N. Muennighoff, C. Ociepa, J. Phang, L. Reynolds, H. Schoelkopf, A. Skowron, L. Sutawika, E. Tang, A. Thite, B. Wang, K. Wang, and A. Zou. A framework for few-shot language model evaluation, 07 2024b.
- Google. Our next-generation model: Gemini 1.5. https://blog.google/technology/ai/google-gemini-next-generation-model-february-2024, 2024.
- Google. Gemini 2.5: Our most intelligent ai model. https://blog.google/technology/google-deepmind/gemini-model-thinking-updates-march-2025/, 2025.
- G. T. Google. Gemini: A family of highly capable multimodal models. *ArXiv e-prints*, arXiv:2312.11805, 2023.
- S. Gugger, L. Debut, T. Wolf, P. Schmid, Z. Mueller, S. Mangrulkar, M. Sun, and B. Bossan. Accelerate: Training and inference at scale made simple, efficient and adaptable., 2022.
- S. Gunasekar, Y. Zhang, J. Aneja, C. C. T. Mendes, A. D. Giorno, S. Gopi, M. Javaheripi, P. Kauffmann, G. de Rosa, O. Saarikivi, A. Salim, S. Shah, H. S. Behl, X. Wang, S. Bubeck, R. Eldan, A. T. Kalai, Y. T. Lee, and Y. Li. Textbooks are all you need. *ArXiv e-prints*, arXiv:2306.11644, 2023.
- D. Guo, Q. Zhu, D. Yang, Z. Xie, K. Dong, W. Zhang, G. Chen, X. Bi, Y. Wu, Y. Li, F. Luo, Y. Xiong, and W. Liang. Deepseek-coder: When the large language model meets programming the rise of code intelligence. *ArXiv e-prints*, arXiv:2401.14196, 2024a.
- H. Guo, Y. Yao, W. Shen, J. Wei, X. Zhang, Z. Wang, and Y. Liu. Human-instruction-free llm self-alignment with limited samples. *ArXiv e-prints*, arXiv:2401.06785, 2024b.
- H. Gupta, K. Scaria, U. Anantheswaran, S. Verma, M. Parmar, S. A. Sawant, C. Baral, and S. Mishra. TarGEN: Targeted data generation with large language models. *ArXiv e-prints*, arXiv:2310.17876, 2023.
- E. Hazan. Introduction to online convex optimization. *Foundations and TrenRds in Optimization*, 2 (3–4):157–325, 2016.
- D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt. Measuring massive multitask language understanding. In *Proceedings of the 9th International Conference on Learning Representations*, 2021a.
- D. Hendrycks, C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, and J. Steinhardt. Measuring mathematical problem solving with the MATH dataset. *ArXiv e-prints*, arXiv:2103.03874, 2021b.
- J. Huang and K. C.-C. Chang. Towards reasoning in large language models: A survey. In *Findings of the 61st Conference of the Association for Computational Linguistics*, pages 1049–1065, 2023.
- B. Hui, J. Yang, Z. Cui, J. Yang, D. Liu, L. Zhang, T. Liu, J. Zhang, B. Yu, K. Lu, K. Dang, Y. Fan, Y. Zhang, A. Yang, R. Men, F. Huang, B. Zheng, Y. Miao, S. Quan, Y. Feng, X. Ren, X. Ren, J. Zhou, and J. Lin. Qwen2.5-Coder technical report. *ArXiv e-prints*, arXiv:2409.12186, 2024.
- A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. de las Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, L. R. Lavaud, M.-A. Lachaux, P. Stock, T. L. Scao, T. Lavril, T. Wang, T. Lacroix, and W. E. Sayed. Mistral 7B. *ArXiv e-prints*, arXiv:2310.06825, 2023.
- J. Jiang, F. Wang, J. Shen, S. Kim, and S. Kim. A survey on large language models for code generation. *ArXiv e-prints*, arXiv:2406.00515, 2024a.
- J.-P. Jiang, Y. Xia, H.-L. Sun, S. Lu, Q.-G. Chen, W. Luo, K. Zhang, D.-C. Zhan, and H.-J. Ye. Multi-modal tabular reasoning with privileged structured information. *ArXiv e-prints*, arXiv:2506.04088, 2025.
- X. Jiang, Y. Ge, Y. Ge, D. Shi, C. Yuan, and Y. Shan. Supervised fine-tuning in turn improves visual foundation models. *ArXiv e-prints*, arXiv:2401.10222, 2024b.

- M. Josifoski, M. Sakota, M. Peyrard, and R. West. Exploiting asymmetry for synthetic training data generation: Synthie and the case of information extraction. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1555–1574, 2023.
- M. Lanctot, V. Zambaldi, A. Gruslys, A. Lazaridou, K. Tuyls, J. Perolat, D. Silver, and T. Graepel. A unified game-theoretic approach to multiagent reinforcement learning. In *Advances in Neural Information Processing Systems 30*, pages 4190–4203, 2017.
- H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning. In *Advances in Neural Information Processing Systems 36*, pages 34892–34916, 2023.
- H. Liu, C. Li, Y. Li, B. Li, Y. Zhang, S. Shen, and Y. J. Lee. LLaVA-NeXT: Improved reasoning, ocr, and world knowledge. https://llava-vl.github.io/blog/2024-01-30-llava-next/, 2024.
- L. Long, R. Wang, R. Xiao, J. Zhao, X. Ding, G. Chen, and H. Wang. On LLMs-driven synthetic data generation, curation, and evaluation: A survey. *ArXiv e-prints*, arXiv:2406.15126, 2024.
- S. Lu, Y. Li, Y. Xia, Y. Hu, S. Zhao, Y. Ma, Z. Wei, Y. Li, L. Duan, J. Zhao, Y. Han, H. Li, W. Chen, J. Tang, C. Hou, Z. Du, T. Zhou, W. Zhang, H. Ding, J. Li, W. Li, G. Hu, Y. Gu, S. Yang, J. Wang, H. Sun, Y. Wang, H. Sun, J. Huang, Y. He, S. Shi, W. Zhang, G. Zheng, J. Jiang, S. Gao, Y.-F. Wu, S. Chen, Y. Chen, Q.-G. Chen, Z. Xu, W. Luo, and K. Zhang. Ovis2.5 technical report. *ArXiv e-prints*, arXiv:2508.11737, 2025.
- Z. Luo, C. Xu, P. Zhao, Q. Sun, X. Geng, W. Hu, C. Tao, J. Ma, Q. Lin, and D. Jiang. Wizard-Coder: Empowering code large language models with evol-instruct. In *Proceedings of the 12th International Conference on Learning Representations*, 2024.
- P. Muller, S. Omidshafiei, M. Rowland, K. Tuyls, J. Perolat, S. Liu, D. Hennes, L. Marris, M. Lanctot, E. Hughes, Z. Wang, G. Lever, N. Heess, T. Graepel, and R. Munos. A generalized training approach for multiagent learning. *ArXiv e-prints*, arXiv:1909.12823, 2019.
- A. Müller. Integral probability metrics and their generating classes of functions. Advances in Applied Probability, 29(2):429–443, 1997.
- OpenAI. GPT-4. https://openai.com/index/gpt-4-research/, 2023.
- OpenAI. Hello GPT-4o. https://openai.com/index/hello-gpt-4o/, 2024.
- OpenAI. Introducing GPT-4.5. https://openai.com/index/introducing-gpt-4-5/, 2025.
- F. Orabona. A modern introduction to online learning. ArXiv e-prints, arXiv:1912.13213v6, 2019.
- L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P. Christiano, J. Leike, and R. Lowe. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems 35*, pages 27730–27744, 2022.
- R. Y. Pang, W. Yuan, H. He, K. Cho, S. Sukhbaatar, and J. Weston. Iterative reasoning preference optimization. In *Advances in Neural Information Processing Systems* 37, 2024.
- B. Peng, C. Li, P. He, M. Galley, and J. Gao. Instruction tuning with GPT-4. *ArXiv e-prints*, arXiv:2304.03277, 2023.
- J. Piaget. The origins of intelligence in children. International Universities Press, 1952.
- Y. Qin, S. Liang, Y. Ye, K. Zhu, L. Yan, Y. Lu, Y. Lin, X. Cong, X. Tang, B. Qian, S. Zhao, L. Hong, R. Tian, R. Xie, J. Zhou, M. Gerstein, D. Li, Z. Liu, and M. Sun. ToolLLM: Facilitating large language models to master 16000+ real-world APIs. *ArXiv e-prints*, arXiv:2307.16789, 2023.
- Z. Qin, Q. Dong, X. Zhang, L. Dong, X. Huang, Z. Yang, M. Khademi, D. Zhang, H. H. Awadalla, Y. R. Fung, W. Chen, M. Cheng, and F. Wei. Scaling laws of synthetic data for language models. *ArXiv e-prints*, arXiv:2503.19551, 2025.

- R. Rafailov, A. Sharma, E. Mitchell, C. D. Manning, S. Ermon, and C. Finn. Direct preference optimization: Your language model is secretly a reward model. In *Advances in Neural Information Processing Systems* 36, pages 53728–53741, 2023.
- D. Rein, B. L. Hou, A. C. Stickland, J. Petty, R. Y. Pang, J. Dirani, J. Michael, and S. R. Bowman. GPQA: A graduate-level google-proof Q&A benchmark. *ArXiv e-prints*, arXiv:2311.12022, 2023.
- M. Ren, B. Cao, H. Lin, C. Liu, X. Han, K. Zeng, G. Wan, X. Cai, and L. Sun. Learning or self-aligning? rethinking instruction fine-tuning. *ArXiv e-prints*, arXiv:2402.18243, 2024.
- C. Rosset, C.-A. Cheng, A. Mitra, M. Santacroce, A. Awadallah, and T. Xie. Direct nash optimization: Teaching language models to self-improve with general preferences. *ArXiv e-prints*, arXiv:2404.03715, 2024.
- K. Sakaguchi, R. L. Bras, C. Bhagavatula, and Y. Choi. WinoGrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- P. Z. Scroccaro, A. S. Kolarijani, and P. M. Esfahani. Adaptive composite online optimization: Predictions in static and dynamic environments. *IEEE Transactions on Automatic Control*, 68(5): 2906–2921, 2023.
- D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. P. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *ArXiv e-prints*, arXiv:1712.01815, 2017a.
- D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, Y. Chen, T. Lillicrap, F. Hui, L. Sifre, G. van den Driessche, T. Graepel, and D. Hassabis. Mastering the game of go without human knowledge. *Nature*, 550(7676):354–359, 2017b.
- Z. Sprague, X. Ye, K. Bostrom, S. Chaudhuri, and G. Durrett. MuSR: Testing the limits of chain-of-thought with multistep soft reasoning. *ArXiv e-prints*, arXiv:2310.16049, 2024.
- H.-L. Sun, Z. Sun, H. Peng, and H.-J. Ye. Mitigating visual forgetting via take-along visual conditioning for multi-modal long cot reasoning. In *Proceedings of the 63rd Conference of the Association for Computational Linguistics*, pages 5158–5171, 2025a.
- H.-L. Sun, D.-W. Zhou, Y. Li, S. Lu, C. Yi, Q.-G. Chen, Z. Xu, W. Luo, K. Zhang, and D.-C. Zhan. Parrot: Multilingual visual instruction tuning. In *Proceedings of the 42nd International Conference on Machine Learning*, 2025b.
- M. Suzgun, N. Scales, N. Schärli, S. Gehrmann, Y. Tay, H. W. Chung, A. Chowdhery, Q. V. Le, E. H. Chi, D. Zhou, and J. Wei. Challenging BIG-Bench tasks and whether chain-of-thought can solve them. *ArXiv e-prints*, arXiv:2210.09261, 2022.
- Z. Tan, D. Li, S. Wang, A. Beigi, B. Jiang, A. Bhattacharjee, M. Karami, J. Li, L. Cheng, and H. Liu. Large language models for data annotation and synthesis: A survey. *ArXiv e-prints*, arXiv:2402.13446, 2024.
- R. Taori, I. Gulrajani, T. Zhang, Y. Dubois, X. Li, C. Guestrin, P. Liang, and T. B. Hashimoto. Stanford Alpaca: An instruction-following LLaMA model, 2023.
- G. Tesauro. Temporal difference learning and TD-Gammon. *Communications of the ACM*, 38(3): 58–68, 1995.
- L. Tunstall, E. Beeching, N. Lambert, N. Rajani, K. Rasul, Y. Belkada, S. Huang, L. von Werra, C. Fourrier, N. Habib, N. Sarrazin, O. Sanseviero, A. M. Rush, and T. Wolf. Zephyr: Direct distillation of LM alignment. *ArXiv e-prints*, arXiv:2310.16944, 2023.
- L. Tunstall, E. Beeching, N. Lambert, N. Rajani, S. Huang, K. Rasul, A. Bartolome, A. M. Rush, and T. Wolf. The Alignment Handbook, 2024. URL https://github.com/huggingface/alignment-handbook.

- S. Ubani, S. O. Polat, and R. Nielsen. ZeroShotDataAug: Generating and augmenting training data with chatgpt. *ArXiv e-prints*, arXiv:2304.14334, 2023.
- K. Wang, J. Zhu, M. Ren, Z. Liu, S. Li, Z. Zhang, C. Zhang, X. Wu, Q. Zhan, Q. Liu, and Y. Wang. A survey on data synthesis and augmentation for large language models. *ArXiv e-prints*, arXiv:2410.12896, 2024a.
- P. Wang, M. Wang, Z. Ma, X. Yang, S. Feng, D. Wang, Y. Zhang, and K. Song. Language models as continuous self-evolving data engineers. *ArXiv e-prints*, arXiv:2412.15151, 2024b.
- Y. Wang, S. Chen, W. Jiang, W. Yang, Y. Wan, and L. Zhang. Online composite optimization between stochastic and adversarial environments. In *In Advances in Neural Information Processing Systems* 37, pages 94808–94850, 2024c.
- Y. Wang, X. Ma, G. Zhang, Y. Ni, A. Chandra, S. Guo, W. Ren, A. Arulraj, X. He, Z. Jiang, T. Li, M. Ku, K. Wang, A. Zhuang, R. Fan, X. Yue, and W. Chen. MMLU-Pro: A more robust and challenging multi-task language understanding benchmark. *ArXiv e-prints*, arXiv:2406.01574, 2024d.
- Y. Wang, W. Yang, W. Jiang, S. Lu, B. Wang, H. Tang, Y. Wan, and L. Zhang. Non-stationary projection-free online learning with dynamic and adaptive regret guarantees. In *Proceedings of the 38th AAAI Conference on Artificial Intelligence*, pages 15671–15679, 2024e.
- Y. Wang, G. Huzhang, Q.-G. Chen, Z. Xu, W. Luo, K. Zhang, and L. Zhang. SPACE: Noise contrastive estimation stabilizes self-play fine-tuning for large language models. In *Advances in Neural Information Processing Systems 38*, 2025.
- L. Wu, Z. Zheng, Z. Qiu, H. Wang, H. Gu, T. Shen, C. Qin, C. Zhu, H. Zhu, Q. Liu, H. Xiong, and E. Chen. A survey on large language models for recommendation. *World Wide Web*, 27(60), 2024a.
- Y. Wu, Z. Sun, H. Yuan, K. Ji, Y. Yang, and Q. Gu. Self-play preference optimization for language model alignment. *ArXiv e-prints*, arXiv:2405.00675, 2024b.
- H. Xin, D. Guo, Z. Shao, Z. Ren, Q. Zhu, B. Liu, C. Ruan, W. Li, and X. Liang. DeepSeek-Prover: Advancing theorem proving in llms through large-scale synthetic data. *ArXiv e-prints*, arXiv:2405.14333, 2024.
- B. Xu, A. Yang, J. Lin, Q. Wang, C. Zhou, Y. Zhang, and Z. Mao. ExpertPrompting: Instructing large language models to be distinguished experts. *ArXiv e-prints*, arXiv:2305.14688, 2023.
- A. Yang, B. Zhang, B. Hui, B. Gao, B. Yu, C. Li, D. Liu, J. Tu, J. Zhou, J. Lin, K. Lu, M. Xue, R. Lin, T. Liu, X. Ren, and Z. Zhang. Qwen2.5-math technical report: Toward mathematical expert model via self-improvement. *ArXiv e-prints*, arXiv:2409.12122v1, 2024.
- H. Yuan, Z. Chen, K. Ji, and Q. Gu. Self-play fine-tuning of diffusion models for text-to-image generation. *ArXiv e-prints*, arXiv:2402.10210, 2024a.
- W. Yuan, R. Y. Pang, K. Cho, X. Li, S. Sukhbaatar, J. Xu, and J. Weston. Self-rewarding language models. In *Proceedings of the 41st International Conference on Machine Learning*, pages 57905–57923, 2024b.
- R. Zellers, A. Holtzman, Y. Bisk, A. Farhadi, and Y. Choi. HellaSwag: Can a machine really finish your sentence? In *Proceedings of the 57th Conference of the Association for Computational Linguistics*, pages 4791–4800, 2019.
- L. Zhang, S. Lu, and Z.-H. Zhou. Adaptive online learning in dynamic environments. In *Advances in Neural Information Processing Systems 31*, pages 1323–1333, 2018.
- Y. Zhang. Training and evaluating language models with template-based data generation. *ArXiv e-prints*, arXiv:2411.18104, 2024.
- Z. Zhao, W. Fan, J. Li, Y. Liu, X. Mei, Y. Wang, Z. Wen, F. Wang, X. Zhao, J. Tang, and Q. Li. Recommender systems in the era of large language models (LLMs). *IEEE Transactions on Knowledge and Data Engineering*, 36:6889–6907, 2024.

- J. Zhou, T. Lu, S. Mishra, S. Brahma, S. Basu, Y. Luan, D. Zhou, and L. Hou. Instruction-following evaluation for large language models. *ArXiv e-prints*, arXiv:2311.07911, 2023.
- Y. Zhu, P. Zhang, E.-U. Haq, P. Hui, and G. Tyson. Can ChatGPT reproduce human-generated labels? a study of social computing tasks. *ArXiv e-prints*, arXiv:2304.10145, 2023.

### **NeurIPS Paper Checklist**

#### 1. Claims

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

Answer: [Yes]

Justification: We have included the main contributions in the abstract and introduction.

Guidelines:

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

#### 2. Limitations

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

Answer: [Yes]

Justification: The limitations have been discussed in the Section 6.

Guidelines:

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

#### 3. Theory assumptions and proofs

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

Answer: [Yes]

Justification: We have specified the assumptions of Theorem 1 and provided the complete proof in the Appendix A.2.

#### Guidelines:

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

### 4. Experimental result reproducibility

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

Answer: [Yes]

Justification: All information about the experiments is provided in the Appendix B.1.

#### Guidelines:

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

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

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

Answer: [No]

Justification: The complete details of the experimental setup are provided in the Section 4 and Appendix B.1, and the pseudocode are included in the Appendix D.

#### Guidelines:

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

#### 6. Experimental setting/details

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

Answer: [Yes]

Justification: The information has been provided in the Appendix B.1.

#### Guidelines:

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

#### 7. Experiment statistical significance

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

Answer: [No]

Justification: Error bars are not reported because it would be too computationally expensive. Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)

- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

#### 8. Experiments compute resources

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

Answer: [Yes]

Justification: The information has been provided in the Appendix B.1.

#### Guidelines:

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

#### 9. Code of ethics

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

Answer: [Yes]

Justification: The authors have diligently adhered to the NeurIPS Code of Ethics in all aspects of this research.

### Guidelines:

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

### 10. Broader impacts

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

Answer: [NA]

Justification: There is no societal impact of this work.

#### Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

#### 11. Safeguards

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

Answer: [NA]

Justification: Our work does not involve releasing data or models that have a high risk for misuse.

#### Guidelines:

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

#### 12. Licenses for existing assets

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

Answer: [Yes]

Justification: We have properly credited the original owners of the assets.

#### Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

• If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

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

Answer: [NA]

Justification: This paper does not release new assets.

#### Guidelines:

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

### 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

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

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

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

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

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

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

### 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: This paper is a methodological research on the self-play fine-tuning of LLMs, and has included the specification about training LLMs in Appendix B.1.

#### Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

#### A Mathematical derivations

In this section, we provide mathematical derivations for the propositions and theorems.

#### A.1 Proof of Proposition 1

In this part, we proceed to derive the optimal solution of the following objective:

$$\pi_{\hat{\theta}} = \operatorname*{argmax}_{\pi_{\theta}} \mathbb{E}_{\mathbf{x} \sim q(\cdot), \mathbf{y}' \sim \pi_{\theta}(\cdot|\mathbf{x})} \left[ c_{t+1}(\mathbf{x}, \mathbf{y}') \right] + \alpha \mathbb{E}_{\mathbf{x} \sim q(\cdot)} \left[ \mathcal{H}(\pi_{\theta}(\cdot|\mathbf{x})) \right]. \tag{10}$$

under a general confidence function  $c_{t+1}$  with the entropy constraint  $\mathcal{H}(\cdot)$ . The proof is inspired by Rafailov et al. [2023], converting (10) into a KL-divergence minimization problem. Specifically, by the definition of entropy, we have:

$$\mathcal{H}(\pi_{\theta}(\cdot|\mathbf{x})) = -\mathbb{E}_{\mathbf{y}' \sim \pi_{\theta}(\cdot|\mathbf{x})} \left[ \log \pi_{\theta}(\mathbf{y}'|\mathbf{x}) \right]. \tag{11}$$

Combining (10) and (11), we have:

$$\underset{\pi_{\theta}}{\operatorname{argmax}} \mathbb{E}_{\mathbf{x} \sim q(\cdot), \mathbf{y}' \sim \pi_{\theta}(\cdot | \mathbf{x})} \left[ c_{t+1}(\mathbf{x}, \mathbf{y}') \right] + \alpha \mathbb{E}_{\mathbf{x} \sim q(\cdot)} \left[ \mathcal{H}(\pi_{\theta}(\cdot | \mathbf{x})) \right] \\
= \underset{\pi_{\theta}}{\operatorname{argmax}} \mathbb{E}_{\mathbf{x} \sim q(\cdot), \mathbf{y}' \sim \pi_{\theta}(\cdot | \mathbf{x})} \left[ c_{t+1}(\mathbf{x}, \mathbf{y}') - \alpha \log \pi_{\theta}(\mathbf{y}' | \mathbf{x}) \right] \\
= \underset{\pi_{\theta}}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x} \sim q(\cdot), \mathbf{y}' \sim \pi_{\theta}(\cdot | \mathbf{x})} \left[ \log \pi_{\theta}(\mathbf{y}' | \mathbf{x}) - \frac{1}{\alpha} c_{t+1}(\mathbf{x}, \mathbf{y}') \right] \\
= \underset{\pi_{\theta}}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x} \sim q(\cdot), \mathbf{y}' \sim \pi_{\theta}(\cdot | \mathbf{x})} \left[ \log \frac{\pi_{\theta}(\mathbf{y}' | \mathbf{x})}{\frac{1}{Z(\mathbf{x})} \exp(\frac{1}{\alpha} c_{t+1}(\mathbf{x}, \mathbf{y}'))} - \log Z(\mathbf{x}) \right]. \tag{12}$$

where  $Z(\mathbf{x}) = \sum_{\mathbf{z} \in \mathcal{Y}} \exp(\frac{1}{\alpha} c_{t+1}(\mathbf{x}, \mathbf{z}))$  denotes the partition function for normalization. Note that  $Z(\mathbf{x})$  is a function of the prompt  $\mathbf{x}$ , and independent of the policy  $\pi_{\theta}$ . We define

$$\pi_{\theta^*}(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp\left(\frac{1}{\alpha}c_{t+1}(\mathbf{x},\mathbf{y})\right),$$

which is a valid policy, since  $\pi_{\theta^*}(\mathbf{y}|\mathbf{x}) \geq 0$  for any  $\mathbf{y}$  and  $\sum_{\mathbf{y}} \pi_{\theta^*}(\mathbf{y}|\mathbf{x}) = 1$ . Next, we can rewrite (12) as:

$$\underset{\pi_{\theta}}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x} \sim q(\cdot), \mathbf{y}' \sim \pi_{\theta}(\cdot | \mathbf{x})} \left[ \log \frac{\pi_{\theta}(\mathbf{y}' | \mathbf{x})}{\pi_{\theta^*}(\mathbf{y}' | \mathbf{x})} - \log Z(\mathbf{x}) \right]$$

$$= \underset{\pi_{\theta}}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x} \sim q(\cdot)} \left[ \mathbb{D}_{KL}(\pi_{\theta}(\cdot | \mathbf{x}) || \pi_{\theta^*}(\cdot | \mathbf{x})) - \log Z(\mathbf{x}) \right],$$

from which we can observe that the optimal policy is exactly  $\pi_{\theta^*}$ . Therefore, the opponent policy is

$$\pi_{\hat{\theta}}(\mathbf{y}|\mathbf{x}) = \pi_{\theta^*}(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{x})} \exp\left(\frac{1}{\alpha}c_{t+1}(\mathbf{x},\mathbf{y})\right).$$

#### A.2 Proof of Theorem 1

In this part, we derive the gradient of (7) with respect to the parameter  $\theta$ . For brevity, we denote

$$u(\pi_{\theta}) = \log \pi_{\theta}(\mathbf{y}|\mathbf{x}) - \log \pi_{\theta}(\mathbf{y}'|\mathbf{x})$$
  
$$v(\pi_{\theta}) = \log \pi_{\theta}(\mathbf{y}'|\mathbf{x}) - \log \pi_{\theta}(\mathbf{y}_{0}|\mathbf{x})$$
(13)

and the expectation operator in this part is taken over the distributions  $\mathbf{x} \sim q(\cdot)$ ,  $\mathbf{y} \sim \pi_{data}(\cdot|\mathbf{x})$ ,  $\mathbf{y}' \sim \pi_{\theta_t}(\cdot|\mathbf{x})$ ,  $\mathbf{y}_0 \sim \pi_{\theta_0}(\cdot|\mathbf{x})$ . By utilizing the chain rule, we have

$$\nabla_{\theta} \mathcal{L}_{\text{T-SPIN}} = \mathbb{E} \left[ \nabla_{\theta} \ell(\alpha u(\pi_{\theta})) + \beta \nabla_{\theta} \ell(\alpha v(\pi_{\theta})) \right]$$

$$= \alpha \mathbb{E} \left[ \ell'(\alpha u(\pi_{\theta})) \nabla_{\theta} u(\pi_{\theta}) \right] + \alpha \beta \mathbb{E} \left[ \ell'(\alpha v(\pi_{\theta})) \nabla_{\theta} v(\pi_{\theta}) \right].$$
(14)

Next, we consider two specific choices of  $\ell(x)$  and its derivative  $\ell'(x)$ . For the logistic loss function  $\ell(x) = -\log \sigma(x)$ , we have  $\ell'(x) = -\sigma(x)$ . Then, combining (13) and (14) delivers

$$\nabla_{\theta} \mathcal{L}_{\text{T-SPIN}} = -\alpha \mathbb{E} \left[ \sigma(\alpha u(\pi_{\theta})) \nabla_{\theta} u(\pi_{\theta}) + \beta \sigma(\alpha v(\pi_{\theta})) \nabla_{\theta} v(\pi_{\theta}) \right]$$

$$= -\alpha \mathbb{E} \left[ \sigma(\alpha u(\pi_{\theta})) \left[ \nabla_{\theta} \log \pi_{\theta}(\mathbf{y}|\mathbf{x}) - \nabla_{\theta} \log \pi_{\theta}(\mathbf{y}'|\mathbf{x}) \right] + \beta \sigma(\alpha v(\pi_{\theta})) \left[ \nabla_{\theta} \log \pi_{\theta}(\mathbf{y}'|\mathbf{x}) - \nabla_{\theta} \log \pi_{\theta}(\mathbf{y}_{\theta}|\mathbf{x}) \right] \right].$$
(15)

Table 3: The number of few-shots and evaluation metrics for each task.

Datasets	GSM8k	MATH	MUSR	MMLU	MMLU-Pro	GPQA	HellaSwag	Winograde	BBH	IFEval
# few-shot	5	4	0	5	5	0	10	5	3	0
Metric	acc	acc	acc_norm	acc	acc	acc_norn	n acc_norm	acc	acc_norm	acc

For the exponential function  $\ell(x) = \exp(-x)$ , we have  $\ell'(x) = -\exp(-x)$ . Then, combining (13) and (14), we can obtain that

$$\nabla_{\theta} \mathcal{L}_{\text{T-SPIN}} = -\alpha \mathbb{E} \left[ \exp(-\alpha u(\pi_{\theta})) \nabla_{\theta} u(\pi_{\theta}) + \beta \exp(-\alpha v(\pi_{\theta})) \nabla_{\theta} v(\pi_{\theta}) \right]$$

$$= -\alpha \mathbb{E} \left[ \exp(-\alpha u(\pi_{\theta})) \left[ \nabla_{\theta} \log \pi_{\theta}(\mathbf{y}|\mathbf{x}) - \nabla_{\theta} \log \pi_{\theta}(\mathbf{y}'|\mathbf{x}) \right]$$

$$+\beta \exp(-\alpha v(\pi_{\theta})) \left[ \nabla_{\theta} \log \pi_{\theta}(\mathbf{y}'|\mathbf{x}) - \nabla_{\theta} \log \pi_{\theta}(\mathbf{y}_{\theta}|\mathbf{x}) \right] \right].$$
(16)

## **B** More details on experiments

In this section, we provide more information about the experiments, including implementation details, task introductions, comparisons between T-SPIN, SPIN and SFT, and more ablation studies.

#### **B.1** Implementation details

We choose the Alignment Handbook [Tunstall et al., 2024] as the codebase for our experiments, and the Accelerate library [Gugger et al., 2022] for distributed training. Following Chen et al. [2024b], we employ the RMSProp [Scroccaro et al., 2023] with no weight decay as the optimizer for training, and set max sequence length to be 2048 and warmup steps to be 10% of the total training steps. All experiments are conducted on  $8\,\mathrm{H}100\,\mathrm{GPU}\mathrm{s}$  with  $80\mathrm{GB}$  memory. On each GPU, we set the batch size as  $8\,\mathrm{and}$  thus the global batch size is  $64.\,\mathrm{mem}$ 

#### **B.2** Tasks

In the following, we provide a detailed description of the tasks used in the experiments:

- **GSM8K** [Cobbe et al., 2021] is a dataset of 8,500 high-quality, linguistically diverse grade school math word problems, divided into 7.5k training problems and 1k test problems. These problems require between 2 and 8 steps to solve, with solutions primarily consisting of a sequence of elementary calculations using basic arithmetic operations (+, -, /, \*). GSM8K is designed to assess the mathematical reasoning abilities of LLMs.
- MATH [Hendrycks et al., 2021b] contains 12, 500 challenging competition mathematics problems. Each problem in MATH has a full step-by-step solution which can be used to teach models to generate answer derivations and explanations.
- MUSR [Sprague et al., 2024] is designed to evaluate multistep, context-dependent reasoning in language models through long-form, algorithmically generated problems approximately 1,000 words in length. The dataset includes tasks such as murder mysteries, spatial reasoning (e.g., object placement), and team assignment optimization, all presented in natural language narratives.
- MMLU [Hendrycks et al., 2021a] evaluates the multitask accuracy of language models on general knowledge and reasoning across 57 diverse subjects, ranging from STEM fields to humanities. In MMLU, each question is accompanied by four choices, and the answer is one of the four choices.
- MMLU-Pro [Wang et al., 2024d] is an enhanced dataset designed to extend the mostly knowledge-driven MMLU benchmark by integrating more challenging, reasoning-focused questions and expanding the choice set from four to ten options.
- **GPQA** [Rein et al., 2023] is a graduate-level, domain-specific benchmark composed of 448 expertly-authored multiple-choice questions in biology, physics, and chemistry. It aims to evaluate the ability of language models to reason about complex scientific concepts.
- **HellaSwag** [Zellers et al., 2019] is a benchmark dataset designed to evaluate the commonsense reasoning and natural language understanding capabilities of machine learning models, particularly in the context of sentence completion tasks.

Table 4: Performance (%) comparisions among three fine-tuning methods, where SPIN and T-SPIN are trained on 50k annotated data from Ultrachat200k, and SFT is trained on the complete dataset. At the average (Avg) column, we also report improvements (or degradations) over the previous iteration. Abbreviations: MMLU-Pro  $\rightarrow$  MMLUP, Winograde  $\rightarrow$  WG, Instruction Following  $\rightarrow$  IF.

Model		Math & Logic			$Multi-Domain\ Knowledge   Commonsense\ Reasoning  $						IF	Avg	
		GSM8K	MATH	MUSR	MMLU	MMLUP	GPQA	HellaSwag	WG	BBH	IFEval		
Mistral-7B		37.68	3.27	41.27	57.86	29.99	29.86	83.24	74.03	44.26	23.63	42.51	
	SFT	43.69	3.08	42.12	58.37	30.18	27.98	83.40	74.11	45.56	33.22	44.17	
SPIN	Iter0	39.59	2.86	41.26	57.97	30.36	29.47	83.40	74.03	44.56	22.54	42.60	
	Iter1	39.67	3.43	39.65	57.63	29.80	28.33	83.66	74.51	42.50	26.04	42.52 <sub>(-0.08)</sub>	
	Iter2	39.52	3.50	39.64	57.49	29.82	28.00	83.69	73.80	41.97	25.56	$42.30_{(-0.22)}$	
	Iter3	37.34	3.71	37.93	57.80	29.64	28.77	83.91	74.19	41.40	18.31	$41.30_{(-1.00)}$	
	Iter4	37.15	4.56	38.59	58.05	29.90	28.73	84.05	74.43	42.88	24.90	$42.32_{(+1.02)}$	
T-SPIN	Iter0	39.67	3.53	40.20	57.74	30.47	28.82	83.98	74.27	44.06	24.36	42.71	
	Iter1	44.03	3.32	38.95	57.48	30.45	29.10	83.41	74.66	44.91	26.88	$43.32_{(+0.61)}$	
	Iter2	44.31	3.43	38.68	57.61	30.52	30.92	83.47	75.30	45.05	26.33	$43.56_{(+0.24)}$	
	Iter3	45.87	3.85	38.59	58.51	31.42	30.83	84.53	73.88	44.97	36.50	$44.90_{(+1.34)}$	
	Iter4	45.92	3.94	38.46	58.55	31.34	31.06	84.60	74.51	44.90	36.90	$45.02_{(+0.12)}$	

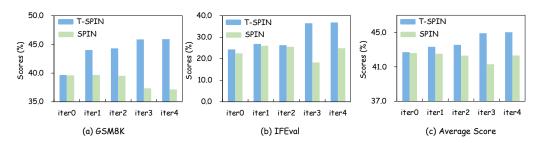


Figure 6: Performance (%) comparisons between T-SPIN and SPIN on two tasks: GMS8K and IFEval over 5 iterations. The average scores over 10 different tasks are also illustrated in the right panel.

- Winograde [Sakaguchi et al., 2021] is a large-scale dataset of 44k problems, inspired by the original WSC design, but adjusted to improve both the scale and the hardness of the dataset.
- **BBH** [Suzgun et al., 2022] is a curated subset of 23 challenging tasks that span algorithmic reasoning, language understanding, world knowledge, and multi-step logic.
- IFEval [Zhou et al., 2023] is a benchmark dataset designed to evaluate the instruction-following capabilities of LLMs. It includes 500 carefully curated prompts annotated with 25 distinct instruction types. The benchmark focuses on a set of "verifiable instructions" such as "write in more than 400 words" and "mention the keyword of AI at least 3 times".

All tasks are evaluated using the Language Model Evaluation Harness framework [Gao et al., 2024b]<sup>3</sup>, with default evaluation settings as specified in Table 3.

#### **B.3** Experimental results on Mistral-7B

In this part, we conduct experiments on Mistral-7B to compare the performances of T-SPIN, SPIN and SFT. The results over 10 different tasks are displayed in Table 4. From the table, we first observe that T-SPIN achieves a notable performance improvement over iterations, and outperforms SFT at the iteration 3. This indicates that with less annotated data, T-SPIN can achieve better performances than SFT that leverages more annotated samples, highlighting the effectiveness of T-SPIN in handling the annotation shortage issue. Second, compared with SPIN, T-SPIN exhibits better performances and stable improvements. Specifically, SPIN suffers from performance degradations at the iterations

<sup>&</sup>lt;sup>3</sup>https://github.com/EleutherAI/lm-evaluation-harness

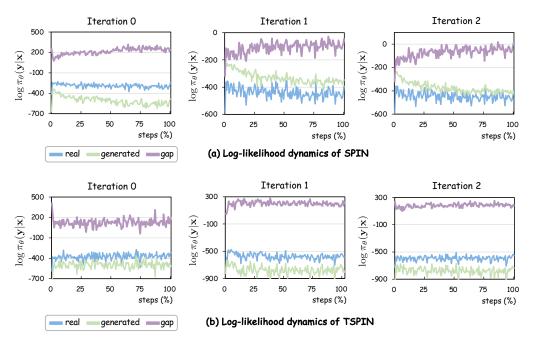


Figure 7: The log-likelihoods of annotated samples and synthetic samples during the model training of T-SPIN and SPIN at different iterations.

1, 2 and 3, while T-SPIN maintains a stable refinement over all iterations. To further illustrate the instability of SPIN, we plot the performances on GSM8K, IFEval and average score in Figure 2. As can be seen, at iteration 3, SPIN exhibits performance degradation compared to iteration 2. In contrast, T-SPIN maintains steady improvements throughout the iterations.

#### B.4 Comparisions on log-likelihoods

In Section 4.2, we have reported the log-likelihood dynamics of two self-play methods at iteration 1. In the following, we examine the trend of log-likelihood changes across different iterations. Specifically, we track the log-likelihoods of annotated samples  $\log \pi_{\theta}(\mathbf{y}|\mathbf{x})$  and synthetic samples  $\log \pi_{\theta}(\mathbf{y}'|\mathbf{x})$  during policy training at iterations 0, 1 and 2, and calculate the gap between them, defined as  $\log \pi_{\theta}(\mathbf{y}|\mathbf{x}) - \log \pi_{\theta}(\mathbf{y}'|\mathbf{x})$ . The results are presented in Figure 7. It is observed that at iteration 0, for both methods, the log-likelihoods of annotated samples are higher than that of synthetic samples, resulting in positive gaps between them. However, at iterations 1 and 2, SPIN demonstrates an unexpected pattern: the log-likelihoods of annotated samples fall below those of synthetic samples, resulting in negative gaps. This reversal stems from the misalignment between the reward function employed in SPIN and the log-likelihood. Specifically, the objective function (3) of SPIN maximizes the reward gap between annotated and synthetic samples. Due to the incorporation of a reference policy, the reward is no longer equivalent to the log-likelihood. As a result, optimizing for high reward during training does not guarantee an increase in log-likelihood. Importantly, since log-likelihood underlies the generation behavior of LLM and reflects its preference for certain responses, it is undesirable that SPIN assigns higher log-likelihoods to synthetic samples than to annotated ones. In contrast, T-SPIN exhibits consistent behavior in terms of log-likelihood: across all iterations, the log-likelihoods of annotated samples remain higher than those of synthetic samples. This desirable property stems from T-SPIN's elimination of the reference policy in its reward function, which ensures alignment between the reward and log-likelihood.

### C Case study

In this section, we provide an example including annotated responses and synthetic responses generated at each iteration, to illustrate the progressive evolution of the model trained by T-SPIN.

Table 5: Examples of the model fine-tuned by T-SPIN at each iteration.

${\tt Prompt} \ \to \\$	Q: Create a thank you card to an organization that helped
	your community in a time of crisis.

#### Annotated ightarrow Dear [Organization Name],

I am writing this letter to express my sincere gratitude to your organization for the exceptional help you provided to our community during a time of crisis. The assistance and support that you and your team showed to us when we were going through difficult times are immeasurable, and we will forever be thankful to you.

Your prompt response, the dedication of your volunteers, and the resources that you made available to us, made a difference in our community. Because of your efforts, families who were affected by [the crisis] found a place to stay, food to eat, and support to help them get back on their feet.

It is remarkable to see how your organization came together and worked tirelessly to provide relief to those in need. Your kindness and empathy for our community are truly appreciated, and it has left a lasting impact on our residents.

Thank you once again for your unwavering support during that trying time. Your actions demonstrate the true meaning of community spirit and the good that can come from working together. Your kindness will always be remembered and cherished.

Sincerely,

### [Your Name]

#### Proto-synthetic ightarrow

1. What is the name of the organization? 2. What is the organization's mission? 3. What is the organization's history? 4. What is the organization's impact on the community? 5. What is the organization's impact on the world? 6. What is the organization's impact on you? 7. What is the organization's impact on your family? 8. What is the organization's impact on your friends? 9. What is the organization's impact on your community? 10. What is the organization's impact on the world?

#### T-SPIN Iter 0 ightarrow

1. I would like to thank the American Red Cross for their help in the aftermath of Hurricane Katrina. The American Red Cross provided shelter, food, and medical care to thousands of people who were displaced by the hurricane. They also provided emotional support to those who had lost loved ones or their homes. The American Red Cross is a vital organization that helps people in times of crisis, and I am grateful for their work. 2. I would also like to thank the Salvation Army for their assistance during Hurricane Katrina. The Salvation Army provided food, water, and shelter to those in need, and they also provided emotional support to those who were struggling. The Salvation Army is a compassionate organization that helps people in times of crisis, and I am grateful for their work. 3. Finally, I would like to thank the United Way for their support during Hurricane Katrina. The United Way provided financial assistance to those who had lost their homes or their jobs, and they also provided emotional support to those who were struggling. The United Way is an organization that helps people in times of crisis, and I am grateful for their work. Thank you to all of these organizations for their help during Hurricane Katrina

- T-SPIN Iter 1  $\to$  1. What is the name of the organization you are thanking? 2. What crisis did this organization help your community with? 3. How did this organization help your community?
  - 4. What did you learn from this experience? 5. How can you help your community in the future?

## Extended Response:

Dear American Red Cross,

I am writing to express my deepest gratitude for the incredible work you did during the recent hurricane that hit our community. Your organization was instrumental in providing relief and support to those affected by the storm. The hurricane caused widespread damage to our area, leaving many people without homes, food, and water. The American Red Cross was there to provide assistance to those in need. Your volunteers worked tirelessly to set up shelters, distribute food and water, and provide emotional support to those who had lost everything.

Your organization's efforts did not go unnoticed. The community was grateful for the help and support you provided during this difficult time. Your actions demonstrated the true spirit of compassion and generosity that is at the heart of the American Red Cross.

Through this experience, I learned the importance of coming together as a community to help those

#### T-SPIN Iter 2 ightarrow

Dear [Organization Name],

I am writing to express my deepest gratitude for the incredible work you did during the recent crisis that affected our community. Your unwavering dedication and commitment to helping those in need truly made a difference in the lives of many.

Your organization's quick response and willingness to go above and beyond to provide assistance to those affected by the crisis was truly inspiring. Your team's tireless efforts to provide food, shelter, and other essential resources to those in need were a true testament to your organization's values and commitment to serving the community.

Your organization's efforts did not go unnoticed, and I am sure that many in our community feel the same way. Your selfless acts of kindness and generosity have touched the hearts of many and have made a lasting impact on our community.

Once again, thank you for your incredible work during this difficult time. Your organization's efforts have not gone unnoticed and are truly appreciated.

Sincerely,

### [Your Name]

T-SPIN Iter 3  $\rightarrow$ Dear [Organization Name],

> I am writing to express my deepest gratitude for the incredible work you did during the recent crisis that affected our community. Your unwavering dedication and commitment to helping those in need truly made a difference in the lives of so many people.

Your organization's quick response and willingness to go above and beyond to provide assistance during a time of crisis is a testament to your selflessness and compassion. Your efforts have not gone unnoticed and are greatly appreciated.

I am proud to be a part of a community that has such an amazing organization like yours. Your tireless work has not only helped those in need but has also brought our community closer together.

Thank you again for all that you do. Your kindness and generosity will never be forgotten.

Sincerely,
[Your Name]

T-SPIN Iter  $4 \rightarrow$  Dear [Organization Name], We, the residents of [ins

We, the residents of [insert community name], would like to extend our deepest gratitude for your unwavering support during the recent natural disaster that struck our town. Your swift response and tireless efforts to provide aid and relief to those affected by the [insert type of disaster] were truly remarkable. Your team of volunteers worked tirelessly to ensure that our community had access to essential resources such as food, water, and shelter. Your kindness and generosity did not go unnoticed and we are forever grateful for the support you provided during this difficult time. You have restored our faith in humanity and have shown us that there is still good in this world. Once again, thank you from the bottom of our hearts for all that you have done for us. You are truly heroes and we are blessed to have had you by our side during this crisis. Sincerely, [Insert community name]

### D Pseudocode

In the following, we provide the pytorch pseudocode for calculating the loss function (7) of T-SPIN.

```
import torch.nn.functional as F
def tpsin_loss(alpha, beta, policy_real_logps, policy_generated_logps,
   policy_proto_logps):
   Compute the loss for T-SPIN (Triplet-based Self-Play Fine-Tuning).
   Args:
       alpha: the regularization coefficient
       beta: the trade-off parameter
       policy_real_logps: Policy's log probs for real samples
       policy_generated_logps: Policy's log probs for synthetic samples
       policy_proto_logps: Policy's log probs for proto-synthetic samples
   Returns:
       torch.Tensor: scalar loss value (mean over batch)
   # Compute advantage terms
   current_advantage = policy_real_logps - policy_generated_logps
   history_advantage = policy_generated_logps - policy_proto_logps
   # Scale advantages to get rewards
   current_rewards = alpha * current_advantage
   history_rewards = alpha * history_advantage
   # Convert rewards to losses using log-sigmoid
   current_loss = -F.logsigmoid(current_rewards)
   history_loss = -F.logsigmoid(history_rewards)
   # Combine loss components
   losses = current_loss + beta * history_loss
   return losses.mean()
```