

Sample, Predict, then Proceed: Self-Verification Sampling for Tool Use of LLMs

Anonymous Author(s)

Affiliation

Address

email

Abstract

Tool use in stateful environments presents unique challenges for large language models (LLMs), where existing test-time compute strategies relying on repeated trials in the environment are impractical. We propose dynamics modelling (DyMo), a method that augments LLMs with a state prediction capability alongside function calling during post-training. This enables LLMs to predict the future states of their actions through an internal environment model. On the Berkeley Function Calling Leaderboard V2, DyMo improves success rates and significantly reduces hallucinations. We further integrate the internal environment model into self-verification sampling (SVS), and show that this substantially improves pass^k over number of trials k , and allows the model to refuse unreliable outputs. Together, DyMo and SVS greatly enhance the effectiveness and reliability of LLMs for tool use. We believe this work charts a path towards scalable planning RL methods for LLM inference without repeatedly querying the oracle environment.

1 Introduction

Large language models (LLMs) have demonstrated remarkable performance in a wide range of applications [1–6]. In addition to conventional natural language tasks, recent advances have shown that LLMs also achieve breakthrough performance in formal language tasks, notably code generation [7–9] and tool use [10–12]. Recent work has shown that scaling the test-time compute can further improve the performance of LLMs on complex tasks such as mathematical reasoning [13–17]. To achieve better performance by scaling up test-time compute, existing methods assume that a verifier, e.g. a process reward model (PRM) or an outcome reward model (ORM), can be queried multiple times during inference [11, 13, 16, 14].

However, many real-world applications may not rely on a verifier to improve test-time sampling, especially when the LLM interacts with the world as in Agentic scenarios. One may not execute k payments and be satisfied that one of the payments is correct among the k ones, whereas one may

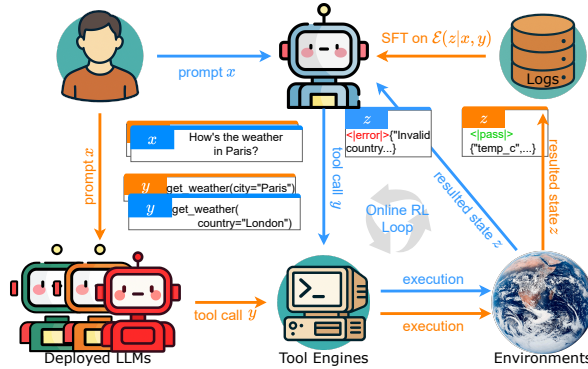


Figure 1: The proposed dynamics modelling (DyMo) that trains an LLM to predict the resulted states of the environment after tools execute function calls via either SFT (orange arrows) on run-logs or over online RL loops (blue arrows).

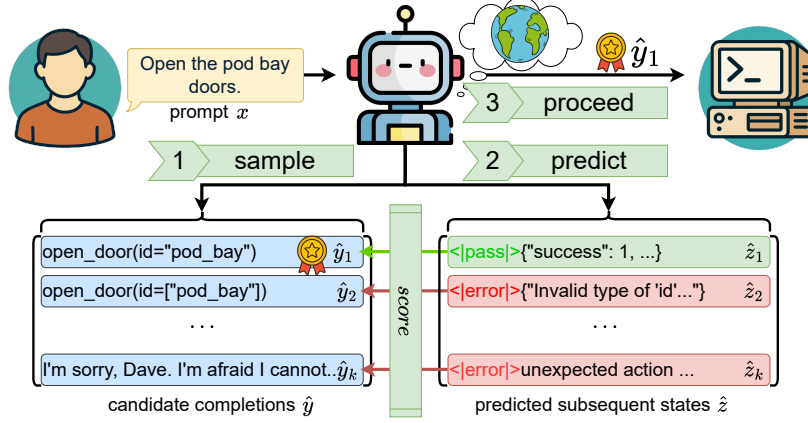


Figure 2: The proposed self-verification sampling (SVS) strategy for test-time compute scaling. For a given user prompt x , the model: 1) generates k candidate completions \hat{y} ; 2) predicts the subsequent states \hat{z} of the k candidates; 3) selects a completion to output by a specified scoring function $score$.

37 verify k times if a mathematical solution is correct without ramifications. In this spirit, we consider
 38 tool-use tasks, where the agent must execute a single trajectory. In particular, such constraints are
 39 often inherent to the **statefulness** of environments, i.e., the environment status states after executing
 40 an action and cannot be easily reverted - the bank account is reduced after a payment!

41 Inspired by the Generative Verifier [18] (GenRM), which formulates the reward function as a next
 42 token prediction task, as illustrated in Figure 1, we propose **dynamic modelling (DyMo)** to fine-tune
 43 LLMs to generate not only the functions calls for a given user prompt, but also the *subsequent states*
 44 of tool engines after executing the generated function calls. This has two advantages: 1) at training
 45 time, this state prediction provides an additional training signal; 2) at test time, this state prediction
 46 can be used in the decision-making process to execute the roll-out, similar to one-step planning
 47 methods.

48 We first explore the impact of DyMo into both the supervised fine-tuning (SFT) and RL stages in
 49 LLM post-training [3–5], and investigate its effectiveness. Our results on the Berkeley Function
 50 Calling Leaderboard [19] V2 (BFCL-V2) show that DyMo alleviates the hallucination problem of
 51 the SFTed model, and improves the success rate of the RLed models. Incorporating DyMo, our
 52 results suggest that an 8B model, when given access to the environment during training, can match
 53 and occasionally surpass the performance of GPT-4o on BFCL-V2.

54 Second, we explore the planning capabilities of DyMo through **self-verification sampling (SVS)**
 55 strategy [20] at test time. Specifically, the models generate k tool calls for a given user prompt, predict
 56 the respective states resulting from those actions, and proceed with the most promising trajectory
 57 based on a ranking mechanism: *sample*, *predict*, then *proceed*. Our experiments demonstrate that
 58 (i) increasing the number of trajectories keeps increasing the LLMs score, (ii) the outcome of the
 59 state prediction can be used to select a successful trajectory *without* access to the oracle environment,
 60 thereby offering a novel schema for scaling test-time compute in stateful environments. Furthermore,
 61 SVS enables models to effectively “*refuse*” requests that exceed their capabilities based on their state
 62 prediction, substantially improving the precision of the final output. We interpret this precision as
 63 “*reliability*”, as it represents the proportion of outputs verified as correct by the oracle environment.

64 In summary, the proposed DyMo method coupled with the SVS strategy significantly enhances the
 65 success rate and reliability of LLMs in tool use tasks.

66 2 Background

67 **Tool Use by LLMs:** Recent works have demonstrated the capability of LLMs to achieve notable
 68 performance in API usage through supervised fine-tuning (SFT) using demonstrations provided either
 69 by human experts or generated by advanced models such as GPT-4 [21–23]. This capability positions
 70 LLMs as back-ends for agents interacting with environments consisting of various tools [19, 24, 25]
 71 and simulated user interactions [26]. However, existing approaches are mainly based on imitation
 72 learning for training [21–23], while the evaluation relies on interactions between environments

and LLMs [26, 19]. Similar to some recent works [27, 28], we focus on learning directly through interacting with the environments, as detailed in Section 3.2.

Reinforcement Learning for Fine-tuning LLMs: Existing RL methodologies for fine-tuning LLMs primarily address alignment tasks [29–31] or reasoning-oriented tasks, such as mathematics and programming challenges [32]. Nonetheless, we posit that RL techniques can effectively extend to tool use scenarios, especially when scaling the quantity of generated tool interactions, given that LLMs have already achieved promising performance in real-world tool use tasks [33, 25, 26]. Furthermore, recent studies indicate a substantial performance gap between online/on-policy RL methods and their offline/off-policy counterparts [34–37]. Although rigorous online interactions can be traded for enhancing wall-clock efficiency, strictly online RL methods still represent an optimal Pareto frontier [37, 38]. Hence, to fully harness the capabilities of RL in tool use contexts, our experiment setup is strictly online and on-policy in this work. Additionally, our method enables models to do one-step planning based their internal learnt environment model during inference time, as illustrated in Section 3.3.

Test-time Compute Scaling: It is well-established that LLMs enhance their performance on logical reasoning tasks by generating extended responses that include explicit intermediate reasoning steps [39]. Further research highlights the importance of explicitly learning these intermediate reasoning stages guided by Policy Reward Models (PRMs) to achieve superior outcomes [14]. While scaling test-time computes by lengthening generated completions has proven beneficial [40, 16], environmental interactions remain critical for achieving optimal results in agent-based tasks [11]. Recent advances also investigate multiple self-rewarding [17], or self-verification [20] steps to scale test-time compute in mathematical reasoning contexts. Unlike these works which query the environment multiple times during inference, we propose to utilise the internal environment models of LLMs to increase the number of completions for scaling test-time computes, as introduced in Section 3.3.

3 Methodology

3.1 Formulation

We used pre-trained Transformer [41] models π_θ parameterised by θ that predict tokens in an autoregressive manner. After post-training by SFT and RL and given a user prompt x , the models can then generate completions/responses y from the distribution $y \sim \pi_\theta(\cdot|x)$. Since we focus on the tool use scenario in this work, we assume the user prompts x are all about requesting function calls, whereas the completions y can be either natural languages or formatted formal languages. For completions that call functions, they will then be passed to the environment \mathcal{E} to execute, and the resulted state is $z = \mathcal{E}(x, y)$ whose complete set is \mathbb{Z} . Note that using no-tool environment is sometimes available in some experiments measuring hallucination.

Following RL terminology, we refer to x as the input *state*, y as the generated *action* from the model π_θ , and z as the resulted next state. The transition dynamics are specified by the environment \mathcal{E} , and the reward function $r : \mathbb{Z} \mapsto [0, 1]$ assigns a binary score to a pair (x, y) according to their resultant state z . From this RL perspective, our model π_θ can:

- generate a tool call (action) y given a user prompt (state) x as input state, i.e. $y \sim \pi_\theta(\cdot|x)$;
- predict next-state z given a user prompt x and a tool call y , i.e. $z \sim \pi_\theta(\cdot|x, y)$.

3.2 DyMo: Dynamics Modelling

The learning objective of the proposed DyMo is not only the tool use function but the environment function \mathcal{E} . As illustrated below, we introduce the DyMo into both the SFT and RL stages.

3.2.1 Dynamics Modelling by Supervised Fine-tuning

During the SFT stage, we construct two distinct datasets — one for the tool use function and one for the environment function — which are described in detail below.

For the tool use function, we train the model π_θ on a dataset of function calls represented by function call (fc) pairs in the form $\langle \text{prompt}, \text{completion} \rangle$, i.e. $\mathbb{D}_{\text{fc}} = \{(x_i, y_i)\}_{i=1}^{N_{\text{fc}}}$. To train the model on these pairs, we minimise the cross-entropy loss [42] of the model’s completion prediction distribution $\pi_\theta(\cdot|x)$ over them:

$$\mathcal{L}_{\text{FC}}(\mathbb{D}_{\text{fc}}; \theta) = - \sum_{i=1}^{N_{\text{fc}}} \sum_{t=1}^{T_{y_i}} \log \pi_{\theta}(y_{i,t} | \mathbf{x}_i, \mathbf{y}_{i,<t}) \quad (1)$$

where $y_{i,t}$ is the t -th element in the target completion \mathbf{y}_i , $\mathbf{y}_{i,<t}$ represents the partial target sequence preceding $y_{i,t}$, and T_{y_i} is the length of \mathbf{y}_i .

Regarding the environment function, we represent it by a dataset of state prediction (sp) triplets in the form $\langle \text{prompt}, \text{completion}, \text{result} \rangle$, i.e. $\mathbb{D}_{\text{sp}} = \{(\mathbf{x}_i, \mathbf{y}_i, \mathbf{z}_i)\}_{i=1}^{N_{\text{sp}}}$. Such data can be gathered and curated from the accumulated *run logs* of the target environment function \mathcal{E} , which we argue is an under-explored source for data scaling. Similar to the tool use function, we minimise the cross-entropy loss of the model’s state prediction distribution $\pi_{\theta}(\cdot | \mathbf{x}, \mathbf{y})$ over these triplets:

$$\mathcal{L}_{\text{SP}}(\mathbb{D}_{\text{sp}}; \theta) = - \sum_{i=1}^{N_{\text{sp}}} \sum_{t=1}^{T_{z_i}} \log \pi_{\theta}(z_{i,t} | \mathbf{x}_i, \mathbf{y}_i, \mathbf{z}_{i,<t}) \quad (2)$$

where the indices i and t follow the same meanings as in Equation 1, and T_{z_i} is the length of \mathbf{z}_i .

3.2.2 Dynamics Modelling over Online Reinforcement Learning

In addition to the SFT stage, DyMo can also be incorporated into the RL fine-tuning of LLMs.

Starting from a prompt set $\mathbb{D}_{\text{rl}} = \{\mathbf{x}_i\}_{i=1}^{N_{\text{rl}}}$, we first sample *two* completions from the model, i.e. $\hat{\mathbf{y}}_i^1, \hat{\mathbf{y}}_i^2 \sim \pi_{\theta}(\cdot | \mathbf{x}_i)$. The two completions along with the prompt \mathbf{x}_i are then passed as inputs to the environment function to get the next states, i.e. $\mathbf{z}_i^1 = \mathcal{E}(\mathbf{x}_i, \hat{\mathbf{y}}_i^1)$ and $\mathbf{z}_i^2 = \mathcal{E}(\mathbf{x}_i, \hat{\mathbf{y}}_i^2)$. Binary scores are then assigned to the $\langle \text{prompt}, \text{completion} \rangle$ pairs by the reward function r , i.e. $r_1 = r(\mathbf{x}_i, \hat{\mathbf{y}}_i^1)$ and $r_2 = r(\mathbf{x}_i, \hat{\mathbf{y}}_i^2)$. Subsequently, we sample predicted next states $\hat{\mathbf{z}}_i^1$ and $\hat{\mathbf{z}}_i^2$ from the model, i.e. $\hat{\mathbf{z}}_i^1 \sim \pi_{\theta}(\cdot | \mathbf{x}_i, \hat{\mathbf{y}}_i^1)$ and $\hat{\mathbf{z}}_i^2 \sim \pi_{\theta}(\cdot | \mathbf{x}_i, \hat{\mathbf{y}}_i^2)$, to track the state prediction performance. Per RL training step, we update the parameter θ of the model π_{θ} to simultaneously minimise the online two-sample REINFORCE Leave-One-Out (RLOO) loss [43, 44] given in Equation 3, and the cross-entropy sample loss given in Equation 4.

$$\mathcal{L}_{\text{RLOO}}(\mathbb{D}_{\text{rl}}; \theta) = - \sum_{i=1}^{N_{\text{rl}}} \left[\left(r_{\beta/2}^{\pi_{\theta}}(\mathbf{x}_i, \hat{\mathbf{y}}_i^1) - r_{\beta/2}^{\pi_{\theta}}(\mathbf{x}_i, \hat{\mathbf{y}}_i^2) \right) \log(\pi_{\theta}(\hat{\mathbf{y}}_i^1 | \mathbf{x}_i)) \right. \\ \left. + \left(r_{\beta/2}^{\pi_{\theta}}(\mathbf{x}_i, \hat{\mathbf{y}}_i^2) - r_{\beta/2}^{\pi_{\theta}}(\mathbf{x}_i, \hat{\mathbf{y}}_i^1) \right) \log(\pi_{\theta}(\hat{\mathbf{y}}_i^2 | \mathbf{x}_i)) \right] \quad (3)$$

$$\mathcal{L}_{\text{DM}}(\mathbb{D}_{\text{rl}}; \theta) = - \sum_{i=1}^{N_{\text{rl}}} \sum_{j=1}^2 \sum_{t=1}^{T_{z_i^j}} \log \pi_{\theta}(z_{i,t}^j | \mathbf{x}_i, \hat{\mathbf{y}}_i^j, \mathbf{z}_{i,<t}^j) \quad (4)$$

where θ_0 is the detached initial parameter in the RL stage, β is a constant hyperparameter, and $r_{\beta/2}^{\pi_{\theta}}(\mathbf{y}_i)$ is the regularised reward defined as $r_j - \frac{\beta}{2} \log \frac{\pi_{\theta}(\mathbf{y}_i^j | \mathbf{x}_i)}{\pi_{\theta_0}(\mathbf{y}_i^j | \mathbf{x}_i)}$ for $j \in \{1, 2\}$.

3.3 Self-Verification Sampling by Internal Environment Model

After undergoing DyMo in both the SFT and RL phases, our model π_{θ} is capable of both generating tool calls and predicting the subsequent states after executing them. Leveraging this capability, we propose to query the internal environment model of π_{θ} multiple times to do Self-Verification Sampling (SVS), as illustrated in Algorithm 1. Given multiple completions per prompt, SVS selects a single output based on a specified scoring function *score* and the internal environment model of π_{θ} . Notably, unlike existing ap-

Algorithm 1: Self-verification sampling (SVS)

Input: x , number of candidate completions k

Output: a completion $\hat{\mathbf{y}}$

Given: a pre-specified scoring function *score*

for $i \leftarrow 1$ **to** k **do**

$\hat{\mathbf{y}}_i \sim \pi_{\theta}(\cdot | \mathbf{x});$

$\hat{\mathbf{z}}_i \sim \pi_{\theta}(\cdot | \mathbf{x}, \hat{\mathbf{y}}_i);$

end

$j \leftarrow \text{score}(\pi_{\theta}(\hat{\mathbf{z}}_1 | \mathbf{x}_1, \hat{\mathbf{y}}_1), \dots, \pi_{\theta}(\hat{\mathbf{z}}_k | \mathbf{x}_k, \hat{\mathbf{y}}_k));$

$\hat{\mathbf{y}} \leftarrow \hat{\mathbf{y}}_j;$

proaches, SVS scales test-time compute *without* querying the oracle environment \mathcal{E} . This approach is reminiscent of the Best-of- N search strategy described in [16], but avoids querying the environment \mathcal{E} multiple times, thereby preventing unintended state changes caused by repeated trials. In addition, SVS aligns with the notion of mental simulation in decision-making, a concept explored in cognitive science [45], thereby establishing a conceptual bridge between research in RL and cognitive science.

4 Experiments

4.1 Setup

Environment: We evaluate tool-use performance using the Berkeley Function Calling Leaderboard V2 (BFCL-V2) [19], which offers comprehensive coverage of function call types, diverse tasks, programming languages, and executability, and has been widely adopted in recent works [46, 6]. As our work is the first to investigate LLMs’ ability to model environment dynamics, we begin with *single-turn* interactions to ensure a clean and tractable problem formulation, in a serverised BFCL-V2 environment in order to run online RL training. Regarding the base model, considering the constraints of our computes, we choose Cohere’s R7B, given its leading performance on various agent benchmarks [6].

SFT Data: During the SFT stage, to constitute the function call (fc) SFT dataset \mathbb{D}_{fc} , inspired by [47, 48], we synthesised pairs of $\langle \text{prompt}, \text{completion} \rangle$ following the distribution of BFCL-V2. Regarding the state prediction (sp) SFT dataset \mathbb{D}_{sp} , we first split the state space \mathbb{Z} into two subsets: 1) pass states \mathbb{Z}^+ where the completions successfully passed the check of BFCL-V2 and received a score of 1; 2) error states \mathbb{Z}^- where the completions failed on the BFCL-V2’s check and received a score of 0. Given the format of BFCL-V2’s return messages, we denote the shared prefix of pass states in \mathbb{Z}^+ as z_{pass} , and similarly z_{error} . Note that, under this setup, there exhibits a bijection between the BFCL-V2’s resulted state subspaces $\{\mathbb{Z}^+, \mathbb{Z}^-\}$ and the scores from the reward function $\{0, 1\}$, which we utilise later to truncate the generation when running SVS during inference time. Following this procedure, we constitute \mathbb{D}_{sp} of $\langle \text{prompt}, \text{completion}, \text{result} \rangle$ triplets from our accumulated run-logs of BFCL-V2 tests.

RL Data: In the following RL stage, to maintain the online RL training and validation distributions as independent and identical, we use 80% from the original BFCL-V2 *prompt* set as the training set, and keep the remaining 20% to validate the generalisation performance. Note that we intentionally keep at least 20 test prompts per category in the final validation set, as certain categories contain ≤ 50 samples, thus 20% of them lacks of statistical significance.

SVS Scoring Function: During the inference time, following GenRM [18], we use a scoring function *score* in the following Equation 5 to run SVS illustrated in Algorithm 1:

$$\text{score}(\pi_{\theta}(\cdot|x_1, \hat{y}_1), \dots, \pi_{\theta}(\cdot|x_k, \hat{y}_k)) \triangleq \arg \max_j (\{\pi_{\theta}(z_{\text{pass}}|x_j, \hat{y}_j)\}). \quad (5)$$

Examples of all the above types of data are provided in Appendix A.

4.2 How proficient is the model at dynamics modelling?

Since we partition the state space \mathbb{Z} to \mathbb{Z}^+ and \mathbb{Z}^- , the state prediction task can be framed as a binary classification problem. A model π^{sft} SFTed on the combined data $\mathbb{D}_{fc} \cup \mathbb{D}_{sp}$, achieves a precision of 90.00%, recall of 87.71%, F1-score of 88.84%, and accuracy of 93.62%. Detailed results are provided in Table 3 in Appendix B. Notably, the success rate of this model on BFCL-V2 is only 72.77%, which is significantly lower than its discriminative performance, highlighting the gap between accurate state prediction and successful functions calls. Therefore, a foundation is laid for improving a model’s generative capability by leveraging its discriminative capability [17, 35].

We also track these metrics during online RL training with the DM loss function for π^{sft} . Under this setting, precision, recall, F1-score, and accuracy further improve to 92.55%, 96.28%, 94.34%, and 90.36%, respectively. The corresponding curves are presented in Figure 6 in Appendix B. For comparison, we SFT an additional model, ϕ^{sft} , on the function call dataset \mathbb{D}_{fc} only. As a result, it learns to roll out states solely through the DM loss during online RL training. We observe a consistent performance gap between this baseline model ϕ^{sft} and π^{sft} , highlighting the necessity of incorporating \mathbb{D}_{sp} in the SFT stage.

4.3 How does dynamics modelling benefit SFT and RL for tool-use?

During the experiments in Section 4.2, we observe that incorporating the additional state prediction data \mathbb{D}_{sp} also leads to a difference in tool use performance. We compare the the performance of ϕ^{sft} - which can do only tool use - with π^{sft} - which is capable of both using tools and predicting

Model	Method	Overall (UW)	Overall (W)	Rel.	Irrel.	AST	Exec
Baselines			<i># samples</i>	(18)	(1122)	(2501)	
GPT-4o [1]	—	82.38	82.14	83.33	81.31	82.51	—
Command-A [6]	—	80.57	84.14	72.22	86.19	83.30	—
Command-R7B	—	70.50	76.70	55.56	81.02	74.92	—
xLAM-2 [48]	—	72.36	71.69	77.78	64.34	74.95	—
ToolACE-2 [23]	—	81.95	85.49	72.22	90.11	83.51	—
watt-tool [49]	—	82.54	81.76	83.33	83.15	81.13	—
BigAgent [50]	—	82.27	81.50	83.33	82.38	81.10	—
SFT			<i># samples</i>	(18)	(1122)	(2501)	
ϕ^{sft}	\mathbb{D}_{fc} only	66.35	66.50	70.73	58.05	70.26	76.25
π^{sft}	$\mathbb{D}_{\text{fc}} \cup \mathbb{D}_{\text{sp}}$	70.87	73.89	63.41	76.32	72.88	77.53
SFT + RL			<i># samples</i>	(20)	(206)	(457)	
ϕ^{rl}	$\phi^{\text{sft}} \rightarrow \text{RLOO}$	80.31	80.22	75.00	89.81	76.13	96.25
ϕ^{rd}	$\phi^{\text{sft}} \rightarrow \text{RLOO} + \text{DyMo}$	82.13	83.16	75.00	91.75	79.65	97.50
π^{rl}	$\pi^{\text{sft}} \rightarrow \text{RLOO}$	81.23	81.99	75.00	90.00	78.68	96.25
π^{rd}	$\pi^{\text{sft}} \rightarrow \text{RLOO} + \text{DyMo}$	83.62	86.68	75.00	90.29	85.56	96.25
SFT + RL + SVS (with k candidates)			<i># samples</i>	(20)	(206)	(457)	
π^{rd}	$\pi^{\text{rd}} \rightarrow \text{SVS with } k = 1$	85.77	84.26	88.20	85.65	83.46	96.33
π^{rd}	$\pi^{\text{rd}} \rightarrow \text{SVS with } k = 2$	88.20	86.71	91.30	86.86	86.44	96.55
π^{rd}	$\pi^{\text{rd}} \rightarrow \text{SVS with } k = 4$	88.94	87.67	91.80	87.41	87.61	96.45
π^{rd}	$\pi^{\text{rd}} \rightarrow \text{SVS with } k = 8$	89.73	88.18	93.10	88.10	88.00	96.25
π^{rd}	$\pi^{\text{rd}} \rightarrow \text{SVS with } k = 16$	89.90	88.29	93.30	88.38	88.03	96.15
π^{rd}	$\pi^{\text{rd}} \rightarrow \text{SVS with } k = 32$	90.18	88.26	94.10	88.59	87.86	96.25
π^{rd}	$\pi^{\text{rd}} \rightarrow \text{SVS with } k = 64$	90.69	88.43	95.00	89.32	87.75	96.25

Table 1: Comprehensive category-wise performance comparison across baselines, SFT, SFT+RL, and SFT+RL+SVS models, on BFCL-V2. For each section, the number of evaluation examples per column is shown in the second row. (W) indicates metrics weighted by the number of samples, whereas UW indicates unweighted. Missing results are marked as ‘—’. The “Exec” column is provided to show the improvement from RL training on it, but is never counted for the overall performance.

next states - across all categories of BFCL-V2. The results in the “SFT” section of Table 1 show that π^{sft} achieves significant improvements on the “Irrelevance” category where the models are not expected to generate function calls when the available tools cannot satisfy the user request. Since the “Irrelevance” is specifically designed to evaluate hallucination of models [19], these results suggest that incorporating the state prediction task helps mitigate hallucination by LLMs [51].

Similarly, we compare the two models — ϕ^{sft} and π^{sft} — both fine-tuned by online RL with and without our DyMo loss, resulting in four variants: π^{rd} , ϕ^{rd} (with DyMo), and π^{rl} , ϕ^{rl} (without DyMo). Take π^{rd} for example, the model is first SFTed on $\mathbb{D}_{\text{fc}} \cup \mathbb{D}_{\text{sp}}$ (thus notated as π), then further fine-tuned by online RL together with DyMo loss (thus superscripted by “rd”). The “SFT + RL” section of Table 1 shows the success rates of these models across different BFCL-V2 categories. For analytical clarity, we preserve the “Exec” category to show the substantial performance gains over it due to RL training. The results indicate that incorporating the DyMo loss yields a $> 5\%$ improvement in success rate over the AST category, contributing to an overall performance boost.

4.4 How does the RL/SFT models perform when scaling up test-time compute?

Since the results in the “SFT” and “SFT + RL” sections of Table 1 are based on a greedy decoding strategy, we further examine whether and how the on-policy distribution over completions for a given prompt, i.e., $\pi_{\theta}(\cdot|x)$, changes under different training pipelines. We begin by analysing the impact of online RL training, comparing the RL-trained models — π^{rl} and ϕ^{rl} — with their corresponding SFT-only baselines — π^{sft} and ϕ^{sft} — using the number of completions per request as the variable. In Figure 3, we report $\text{pass}@k$ and pass^k [26] as the evaluation metrics.

As shown in Figure 3, online RL significantly improves $\text{pass}@k$ when $k \leq 8$. More importantly, online RL consistently improves pass^k over all k values, as evidenced by the consistent gap between the RL-trained models — π^{rl} and ϕ^{rl} — over their SFT-only counterparts — π^{sft} and ϕ^{sft} . These results suggest that the on-policy distributions induced by the RL models yields a more consistent and reliable function calling performance than the distributions induced by the SFT models.

We also note that our $\text{pass}@k$ curves align with the findings in mathematical reasoning tasks [52], where Yue et al. conclude that “base models can achieve a comparable or even higher $\text{pass}@k$ score compared to their RL counterparts at large k values”. However, in our setup, we found that base

model can hardly match pass@ k or pass k of SFT and RL models, which we argue is due to that correct function calls are sparser to generate.

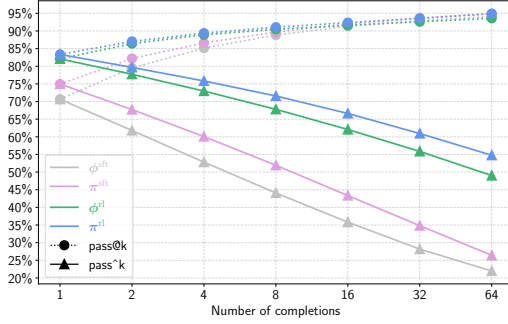


Figure 3: Pass@ k (-) and pass k (..) of RL models - π^{rl} and ϕ^{rl} - vs SFT bases - π^{sft} and ϕ^{sft} on BFCL-V2.

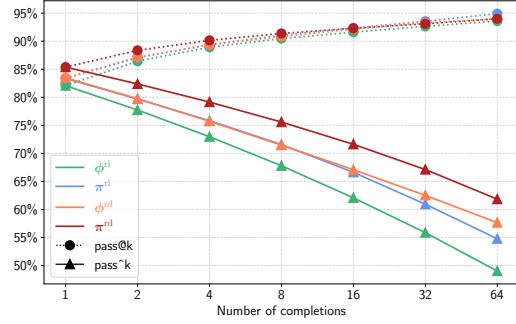


Figure 4: Pass@ k (-) and pass k (..) of models trained by online RL with DyMo loss - π^{rd} and ϕ^{rd} - vs without DyMo loss - π^{rl} and ϕ^{rl} - on BFCL-V2.

4.5 How does dynamics modelling impact the test-time compute scaling of RL models?

Building on the observation that RL models achieve higher success rates over test-time compute scales, we further investigate the impact of incorporating the DyMo loss during online RL training. Similarly, we report pass@ k and pass k for RL models trained with the DyMo loss — π^{rd} and ϕ^{rd} — compared to those trained without it — π^{rl} and ϕ^{rl} .

As shown in Figure 4, adding the DyMo loss during online RL improves pass@ k when $k \leq 8$, while it consistently improves pass k over all numbers of completions per prompt. Note that SVS is not utilised in the experiments so far, thus the improvements are solely due to the DyMo loss. More notably, incorporating the DyMo in both the SFT and RL stages results in π^{rd} , which achieves the highest pass k for all values of k . The consistent gap between pass k curves of π^{rd}/ϕ^{rd} and π^{rl}/ϕ^{rl} also indicate the DyMo loss can help to further improve the consistency and reliability of function calling performance on top of RL. These results demonstrate the effectiveness and benefits of integrating DyMo into both the SFT and RL phases.

4.6 How does self-verification sampling scale over test-time compute?

So far, we have focused primarily on the benefits of incorporating DyMo during model training. However, as introduced in Section 3.3, during inference time, self-verification sampling (SVS) actually unifies the policy (as in model-free RL), the environment model (as in model-based RL), and the value function (under our specific state-space split) into a single LLM. This paradigm enables the model to scale test-time compute by generating more candidate completions per user request **without** querying the oracle environment function \mathcal{E} . To evaluate the effectiveness of SVS, we compare pass k with SVS against pass k without SVS of model π^{rd} . For pass k with SVS, we sample c candidates for each trial and k trials per prompt, thus $k \times c$ candidate completions in total for each prompt. Further, per candidate group for each trial, following GenRM [16], we adopt the scoring function defined in Equation 5 as the metric to select just one output from the c candidates (thus k outputs in the end).

	k	1	2	4	8	16	32
pass k	with SVS	89.02%	87.97%	87.19%	86.14%	84.05%	78.05%
	(c for each trial)	(64)	(32)	(16)	(8)	(4)	(2)
	without SVS	87.68%	82.38%	79.14%	75.58%	71.61%	67.11%

Table 2: Pass k with and without SVS over k trials in the oracle environment. Augmented with SVS, per prompt, we first generate c candidate completions for each trial, then select just one to output by the scoring function defined in Equation 5 for all k trials. Therefore, there are $k \times c$ candidates in total for each prompt, by querying the oracle environment also k times as to pass k without SVS.

As shown in Table 2, SVS achieves improved pass^k over all k values, demonstrating that self-verification enables effective scaling with additional computes. More importantly, the consistent improvement of SVS performance with increasing k highlights our method as a novel test-time compute scaling strategy — one that leverages the model’s internal environment approximation to self-verify and select the most reliable candidate completion. In Section 4.7 and Section 5, we provide further insights about our current SVS setup.

Beyond the above experiment, we also compare $\text{pass}@k$ of the “Best-of- N ” test-time compute scaling strategy with the $\text{pass}@1$ performance of π^{rd} using SVS with k candidate completions per prompt, thus both methods operate under the same inference compute budget. As shown in the results provided in the “SFT + RL + SVS” section of Table 1, increasing number of candidates k in SVS consistently improves the $\text{pass}@1$, which demonstrated the effectiveness of the model’s internal environment model. It is unsurprising to observe that querying the model’s internal environment model is less efficient than accessing the oracle environment function under the same compute budget, and should be seen as an upper-bound. However, we also argue that relying on the oracle may be impractical in many real-world applications involving stateful environments. For example, the model is not expected to place k parallel orders for a single shopping request or to book k tickets on the same flights for a travel planning request.

4.7 What if the model is allowed to refuse?

“I’m sorry, Dave. I’m afraid I can’t do that.” – 2001: A Space Odyssey

As may already be observed, a notable limitation of the scoring function defined in Equation 5 is that the model is still required to output a completion, even in cases where all candidate completions are self-verified as failed trials. That is, our model might roll-out zerror for all generated candidates. In such cases, we argue that it is both reasonable and desirable for the model to “refuse” the request by returning a message that informs the user the query cannot be completed reliably. Formally, we define the revised scoring function as follow:

$$\text{score}(\pi^{\text{rd}}(\hat{z}_1|x_1, \hat{y}_1), \dots, \pi^{\text{rd}}(\hat{z}_k|x_k, \hat{y}_k)) \triangleq \arg \max_j (\{\pi^{\text{rd}}(\hat{z}_j|\hat{y}_j, x_j)\}) \quad (6)$$

s.t. $\text{zpass} \prec \hat{z}_j$ and $\pi^{\text{rd}}(\text{zpass}|\hat{y}_j, x_j) > \tau$

where $\text{zpass} \prec \hat{z}_j$ means that \hat{z}_j starts with zpass as the prefix, and τ is the threshold hyperparameter.

Using Equation 6 as the scoring function, the model π^{rd} classifies a completion \hat{y} as positive if and only if $\pi^{\text{rd}}(\text{zpass}|\hat{y}, x) > \tau$; otherwise, it is classified as negative. By sweeping across a range of thresholds $\tau \in [0.5, 0.99]$, we find that $\tau = 0.92$ offers a favourable trade-off between precision and refusal. Fixing $\tau = 0.92$, we then examine how precision and refuse rate vary with the number of candidate completions k , as shown in Figure 5.

Surprisingly, the model maintains a precision of $\sim 94.5\%$ across values of k , while the refuse rate decreases notably from 23.79% ($k = 1$) to 13.33% ($k = 64$). Since precision reflects the proportion of correct completions among all non-refused outputs, we interpret it as a proxy for the reliability of the model’s responses. Under this view, our results suggest that reliability remains stable as the number of candidates increases, while the refusal rate drops significantly — indicating improved solution coverage without sacrificing correctness. These findings highlight the practical value of combining DyMo with SVS: by generating more candidates, the model achieves higher success rates while maintaining high reliability. We further discuss the broader implications of this observation in Section 5.

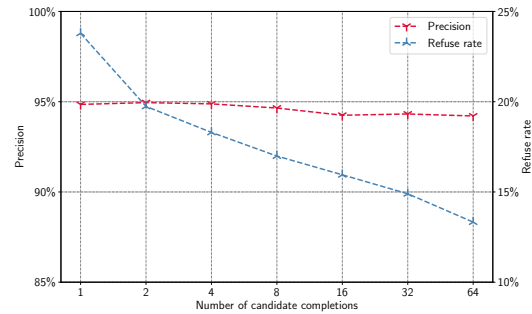


Figure 5: Precision and refuse rate over k for SVS

5 Discussion

Internal environment model by DyMo: Recent advances in RL have shown that incorporating world models can substantially improve performance in complex domains such as board games [53]

and video games [54]. Building on this line of work, our approach takes a further step by unifying the world model and the policy into a single LLM through DyMo, and demonstrates the practical benefits of this unification in tool use scenarios. We also note that similar motivations have emerged in reward modelling, where LLMs are fine-tuned either into stand-alone reward models [18], or into generative models capable of self-rewarding by reasoning over multiple steps in a single completion [17]. Our work extends this broader trend by showing how DyMo can enhance function calling beyond reasoning, particularly in stateful environments.

Low true negative ratio problem of SVS: In our analysis of the results in Section 4.6, we observe that the model π^{rd} exhibits surprisingly low true negative ratio ($< 50\%$ TNR), despite achieving strong precision and recall. As the model’s success rate increases through online RL training, the proportion of correct completions steadily rises. This leads to a highly imbalanced distribution between completions beginning with z_{pass} and z_{error} , thus introduces a bias toward predicting states in \mathbb{Z}^+ . Consequently, we observe that the model tends to “over-refuse” its own completions, i.e. it incorrectly verifies many correct completions as failures via its internal environment model. A straightforward mitigation strategy would be to incorporate additional negative samples from \mathbb{D}_{fc} , thereby exposing the model to a more balanced distribution during DyMo in RL training. Due to time and research scope constraints, we leave this direction for future work. Nonetheless, we argue that our method significantly enhances the reliability of model outputs: higher precision implies that completions self-verified by the model are more likely to be correct. This property is particularly valuable in high-stakes or safety-critical domains, such as healthcare or finance, where even a few incorrect outputs can lead to undesirable or irreversible outcomes.

Test-time compute scaling via DyMo and SVS: As discussed in previous sections, the proposed DyMo and SVS provide a novel strategy for test-time compute scaling. Here, we offer additional reflections from both data-centric and modelling perspectives. First, we highlight that DyMo can benefit from failed completions, since a complete environment model should be capable of handling both successful and failed trajectories. Given the vast amount of run logs accumulated from software systems over decades, we argue that DyMo unlocks a largely under-explored data source: rich, naturally occurring software run logs. In particular, the ability of DyMo to learn from failed completions helps improve the fidelity of the internal environment model — a capability, to the best of our knowledge, not explicitly addressed in prior works from the LLM community. Secondly, from the perspective of world modelling, we hypothesise that programs are often written with an implicit and internal world model of the developers who coded upon assumptions about environment dynamics, constraints, and expected behaviours. These implicit world models are then reflected in the run logs, which can then be captured and fitted by the proposed DyMo method. Through SVS, this learned environment model can be exploited at test time, enabling the model to improve its decision quality without external feedback. While this hypothesis is promising, a deeper exploration of the relationship between program execution and world modelling lies beyond the scope of this work, and we leave it for future investigation.

6 Conclusion

In this work, we investigate the challenge of tool use in stateful environments, where existing test-time compute strategies become impractical due to repeated environment queries. To address this, we propose DyMo, a method that augments LLM fine-tuning with an additional state prediction task during both the SFT and RL stages, enabling a next-state prediction capability of the model. Experiments on the BFCL-V2 benchmark show that incorporating DyMo significantly reduces hallucinations during SFT and improves the success rate over RL training loops. Notably, we also observe that RL models consistently outperform SFT models in mitigating hallucinations. Furthermore, we demonstrate that correct tool calls are retrievable for over 93% of prompts using a parallel Best-of- N decoding strategy, indicating that both SFT and RL models have learned sufficiently expressive on-policy distributions. Building on this insight, we introduce a self-verification sampling (SVS) strategy, which consistently improves pass^k and $\text{pass}@1$ performance by leveraging the model’s internal environment model. Crucially, by allowing the model to refuse uncertain completions, our approach produces more reliable outputs in scenarios where correctness is essential. Overall, our findings highlight a promising direction for extending planning algorithms from the RL community to LLMs in dynamic and stateful environments.

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [2] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [3] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [4] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.
- [5] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- [6] Team Cohere et al. Command a: An enterprise-ready large language model. *arXiv preprint arXiv:2504.00698*, 2025.
- [7] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=VTF8yNQM66>.
- [8] Jonas Gehring, Kunhao Zheng, Jade Copet, Vegard Mella, Quentin Carbonneaux, Taco Cohen, and Gabriel Synnaeve. Rlef: Grounding code llms in execution feedback with reinforcement learning. *arXiv preprint arXiv:2410.02089*, 2024.
- [9] Anthropic. Claude 3.7 sonnet and claude code. <https://www.anthropic.com/news/claude-3-7-sonnet>, 2025. As of April 1st 2025.
- [10] Cohere. Basic usage of tool use (function calling). <https://docs.cohere.com/v2/docs/tool-use-overview#step-3-get-tool-results>, 2024. As of April 1st 2025.
- [11] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023.
- [12] OpenAI. Function calling and other api updates. <https://openai.com/index/function-calling-and-other-api-updates/>, 2023. As of April 1st 2025.
- [13] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652, 2023.
- [14] Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *The Twelfth International Conference on Learning Representations*, 2023.
- [15] Pablo Villalobos and David Atkinson. Trading off compute in training and inference, 2023. URL <https://epochai.org/blog/trading-off-compute-in-training-and-inference>. Accessed, pages 07–03, 2024.
- [16] Charlie Victor Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling LLM test-time compute optimally can be more effective than scaling parameters for reasoning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=4FWAwZtd2n>.

- [17] Wei Xiong, Hanning Zhang, Chenlu Ye, Lichang Chen, Nan Jiang, and Tong Zhang. Self-rewarding correction for mathematical reasoning. *arXiv preprint arXiv:2502.19613*, 2025.
- [18] Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran Kazemi, Aviral Kumar, and Rishabh Agarwal. Generative verifiers: Reward modeling as next-token prediction. *arXiv preprint arXiv:2408.15240*, 2024.
- [19] Fanjia Yan, Huanzhi Mao, Charlie Cheng-Jie Ji, Tianjun Zhang, Shishir G. Patil, Ion Stoica, and Joseph E. Gonzalez. Berkeley function calling leaderboard. https://gorilla.cs.berkeley.edu/blogs/8_berkeley_function_calling_leaderboard.html, 2024.
- [20] Eric Zhao, Pranjal Awasthi, and Sreenivas Gollapudi. Sample, scrutinize and scale: Effective inference-time search by scaling verification. *arXiv preprint arXiv:2502.01839*, 2025.
- [21] Minghao Li, Yingxiu Zhao, Bowen Yu, Feifan Song, Hangyu Li, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. API-bank: A comprehensive benchmark for tool-augmented LLMs. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3102–3116, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.187. URL <https://aclanthology.org/2023.emnlp-main.187/>.
- [22] Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, dahai li, Zhiyuan Liu, and Maosong Sun. ToolLLM: Facilitating large language models to master 16000+ real-world APIs. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=dHng200Jjr>.
- [23] Weiwen Liu, Xu Huang, Xingshan Zeng, xinlong hao, Shuai Yu, Dexun Li, Shuai Wang, Weinan Gan, Zhengying Liu, Yuanqing Yu, Zezhong WANG, Yuxian Wang, Wu Ning, Yutai Hou, Bin Wang, Chuhan Wu, Wang Xinzhi, Yong Liu, Yasheng Wang, Duyu Tang, Dandan Tu, Lifeng Shang, Xin Jiang, Ruiming Tang, Defu Lian, Qun Liu, and Enhong Chen. ToolACE: Winning the points of LLM function calling. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=8EB8k6DdCU>.
- [24] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang. Agentbench: Evaluating LLMs as agents. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=zAdUB0aCTQ>.
- [25] Jiarui Lu, Thomas Holleis, Yizhe Zhang, Bernhard Aumayer, Feng Nan, Felix Bai, Shuang Ma, Shen Ma, Mengyu Li, Guoli Yin, Zirui Wang, and Ruoming Pang. Toolsandbox: A stateful, conversational, interactive evaluation benchmark for llm tool use capabilities, 2024. URL <https://arxiv.org/abs/2408.04682>.
- [26] Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik R Narasimhan. $\{\tau\}$ -bench: A benchmark for Tool-Agent-User interaction in real-world domains. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=roNSXZpUDN>.
- [27] Yuanqing Yu, Zhefan Wang, Weizhi Ma, Zhicheng Guo, Jingtao Zhan, Shuai Wang, Chuhan Wu, Zhiqiang Guo, and Min Zhang. Steptool: A step-grained reinforcement learning framework for tool learning in llms. *arXiv preprint arXiv:2410.07745*, 2024.
- [28] Kevin Chen, Marco Cusumano-Towner, Brody Huval, Aleksei Petrenko, Jackson Hamburger, Vladlen Koltun, and Philipp Krähenbühl. Reinforcement learning for long-horizon interactive llm agents. *arXiv preprint arXiv:2502.01600*, 2025.
- [29] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 3008–3021. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/1f89885d556929e98d3ef9b86448f951-Paper.pdf.

- [30] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=TG8KACxEON>.
- [31] Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, Han Zhao, Yingbo Zhou, Nan Jiang, Doyen Sahoo, Caiming Xiong, and Tong Zhang. Rlhf workflow: From reward modeling to online rlhf. *arXiv preprint arXiv:2405.07863*, 2024.
- [32] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-rl: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [33] Fanjia Yan, Huanzhi Mao, Charlie Cheng-Jie Ji, Tianjun Zhang, Shishir G. Patil, Ion Stoica, and Joseph E. Gonzalez. Berkeley function calling leaderboard. https://gorilla.cs.berkeley.edu/blogs/8_berkeley_function_calling_leaderboard.html, 2024.
- [34] Wei Xiong, Hanze Dong, Chenlu Ye, Ziqi Wang, Han Zhong, Heng Ji, Nan Jiang, and Tong Zhang. Iterative preference learning from human feedback: Bridging theory and practice for rlhf under kl-constraint. *arXiv preprint arXiv:2312.11456*, 2023.
- [35] Shangmin Guo, Biao Zhang, Tianlin Liu, Tianqi Liu, Misha Khalman, Felipe Llinares, Alexandre Rame, Thomas Mesnard, Yao Zhao, Bilal Piot, et al. Direct language model alignment from online ai feedback. *arXiv preprint arXiv:2402.04792*, 2024.
- [36] Yunhao Tang, Daniel Zhaohan Guo, Zeyu Zheng, Daniele Calandriello, Yuan Cao, Eugene Tarassov, Rémi Munos, Bernardo Ávila Pires, Michal Valko, Yong Cheng, et al. Understanding the performance gap between online and offline alignment algorithms. *arXiv preprint arXiv:2405.08448*, 2024.
- [37] Michael Noukhovitch, Shengyi Huang, Sophie Xhonneux, Arian Hosseini, Rishabh Agarwal, and Aaron Courville. Faster, more efficient RLHF through off-policy asynchronous learning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=FhTAG591Ve>.
- [38] Brian R Bartoldson, Siddarth Venkatraman, James Diffenderfer, Moksh Jain, Tal Ben-Nun, Seanie Lee, Minsu Kim, Johan Obando-Ceron, Yoshua Bengio, and Bhavya Kailkhura. Trajectory balance with asynchrony: Decoupling exploration and learning for fast, scalable llm post-training. *arXiv preprint arXiv:2503.18929*, 2025.
- [39] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=_VjQlMeSB_J.
- [40] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- [41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [42] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press Cambridge, 2016.
- [43] Arash Ahmadian, Chris Cremer, Matthias Gallé, Marzieh Fadaee, Julia Kreutzer, Olivier Pietquin, Ahmet Üstün, and Sara Hooker. Back to basics: Revisiting reinforce style optimization for learning from human feedback in llms. *arXiv preprint arXiv:2402.14740*, 2024.

- [44] Yannis Flet-Berliac, Nathan Grinsztajn, Florian Strub, Eugene Choi, Bill Wu, Chris Cremer, Arash Ahmadian, Yash Chandak, Mohammad Gheshlaghi Azar, Olivier Pietquin, and Matthieu Geist. Contrastive policy gradient: Aligning LLMs on sequence-level scores in a supervised-friendly fashion. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 21353–21370, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.1190. URL <https://aclanthology.org/2024.emnlp-main.1190/>.
- [45] Gary Klein and Beth W Crandall. The role of mental simulation in problem solving and decision making. In *Local applications of the ecological approach to human-machine systems*, pages 324–358. CRC Press, 2018.
- [46] Qwen Team. Qwq-32b: Embracing the power of reinforcement learning. <https://qwenlm.github.io/blog/qwq-32b/>, 2024.
- [47] Zuxin Liu, Thai Hoang, Jianguo Zhang, Ming Zhu, Tian Lan, Juntao Tan, Weiran Yao, Zhiwei Liu, Yihao Feng, Rithesh RN, et al. Apigen: Automated pipeline for generating verifiable and diverse function-calling datasets. *Advances in Neural Information Processing Systems*, 37: 54463–54482, 2024.
- [48] Akshara Prabhakar, Zuxin Liu, Weiran Yao, Jianguo Zhang, Ming Zhu, Shiyu Wang, Zhiwei Liu, Tulika Awalgaonkar, Haolin Chen, Thai Hoang, et al. Apigen-mt: Agentic pipeline for multi-turn data generation via simulated agent-human interplay. *arXiv preprint arXiv:2504.03601*, 2025.
- [49] Wentao Shi, Mengqi Yuan, Junkang Wu, Qifan Wang, and Fuli Feng. Direct multi-turn preference optimization for language agents, 2024. URL <https://arxiv.org/abs/2406.14868>.
- [50] BitAgent. Bitagent-8b. <https://huggingface.co/BitAgent/BitAgent-8B>, 2025. Accessed: 2025-05-16.
- [51] Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. On faithfulness and factuality in abstractive summarization. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.173. URL <https://aclanthology.org/2020.acl-main.173/>.
- [52] Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Shiji Song, and Gao Huang. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model? *arXiv preprint arXiv:2504.13837*, 2025.
- [53] John Schultz, Jakub Adamek, Matej Jusup, Marc Lanctot, Michael Kaisers, Sarah Perrin, Daniel Hennes, Jeremy Shar, Cannada Lewis, Anian Ruoss, et al. Mastering board games by external and internal planning with language models. *arXiv preprint arXiv:2412.12119*, 2024.
- [54] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse control tasks through world models. *Nature*, 2025. doi: 10.1038/s41586-025-08744-2. URL <https://doi.org/10.1038/s41586-025-08744-2>.

560 A Examples of Data Used for Training LLMs

561 In this section, we present examples of the datasets curated for supervised fine-tuning (SFT) of large
562 language models (LLMs) on the tool use and state prediction tasks, as described in Section 3.2.1 and
563 Section 4.1.

564 A.1 Example of Function Call Supervised Fine-tuning Dataset \mathbb{D}_{fc}

565 Below, we provide an example of the function call SFT data. The completion shown in the yellow box
566 corresponds to the ground-truth output used for supervised training and is guaranteed to be correct.
567 Importantly, our function call SFT dataset \mathbb{D}_{fc} does not include any data from the original BFCL
568 benchmark; the following example is provided solely for illustrative purposes.

Example 1: Humidity Forecast Query

System Preamble

You are a large language model AI assistant. Your knowledge cutoff date is ...

...

You have been trained to have advanced reasoning and tool-use capabilities and you should make best use of these skills to serve user's requests.

Here is the list of tools that you have available to you. You can ONLY use the tools listed here. When a tool is not listed below, it is NOT available and you should NEVER attempt to use it. Each tool is represented as a JSON object with fields like "name", "description", "parameters" (per JSON Schema), and optionally, "responses" (per JSON Schema).

```
[
{"name": "weather.humidity_forecast", "description": "Retrieve
a humidity forecast for a specific location and time frame.",
"parameters": {"type": "object", "properties": {"location":
{"type": "string", "description": "The city that you want to get the
humidity for."}, "days": {"type": "integer", "description": "Number
of days for the forecast."}, "min_humidity": {"type": "integer",
"description": "Minimum level of humidity (in percentage) to filter
the result. Optional parameter. Default is 0."}}, "required":
["location", "days"]}, "responses": null},
{"name": "get_team_score", "description": "Retrieves the latest game
score, individual player stats, and team stats for a specified sports
team.", "parameters": {"type": "object", "properties": {"team_name":
{"type": "string", "description": "The name of the sports team."},
"league": {"type": "string", "description": "The league that the
team is part of."}, "include_player_stats": {"type": "boolean",
"default": false, "description": "Indicates if individual player
statistics should be included in the result. Default is false."}},
"required": ["team_name", "league"]}, "responses": null}
]
...
```

User Prompt

What is the humidity level in Miami, Florida in the upcoming 7 days?

569

Completion

I'll query the `weather.humidity_forecast` to give user an answer. The call is
`{"weather.humidity_forecast": {"location": ["Miami", "Miami, Florida", "FL"], "days": [7], "min_humidity": ["", 0]}}`

570

571 A.2 Examples of State Prediction Supervised Fine-tuning Dataset \mathbb{D}_{sp}

572 Below, we present an example of the state prediction SFT data. Example 2 corresponds to a previously
573 successful call to the humidity forecast function.

Example 2: A Successful Call to the Humidity Forecast Query

System Preamble

You are a large language model AI assistant. Your knowledge cutoff date is ...
(Identitcal to expmale 1 above, thus we skip the following content to save pages.)

User Prompt

What is the humidity level in Miami, Florida in the upcoming 7 days?

Completion

I'll query the `weather.humidity_forecast` to give user an answer. The call is
`{"weather.humidity_forecast": {"location": ["Miami", "Miami, Florida", "FL"], "days": [7], "min_humidity": ["", 0]}}`

Pass

`<|pass>`
`{ "status": 1, "forecast": [{ "date": "2025-04-11",
"min_humidity": 62, "max_humidity": 78 }, ...] }`

574

575 Example 3 corresponds to a previously failed call to the humidity forecast function and is used as a
576 negative example in the state prediction task.

Example 3: A Failed Call to the Humidity Forecast Query

System Preamble

You are a large language model AI assistant. Your knowledge cutoff date is ...
(Identitcal to expmale 1 above, thus we skip the following content to save pages.)

User Prompt

What is the humidity level in Miami, Florida in the upcoming 7 days?

577

Completion

I'll query the `weather.humidity_forecast` to give user an answer. The call is
`{"weather.humidity_forecast": {"location": ["Miami", "Miami, Florida", "FL"], "days": [-7], "min_humidity": ["", 0]}}`

Error

`<error!>`
`{ "status": 0, "error": "Invalid input to argument 'days', expected a positive integer, but got -7...." }`

578

B Results of Binary Classification

579

580 In Table 3, we present a detailed breakdown of the model’s performance on the binary classification
 581 task formulated in Section 4.2. The table includes the confusion matrix, with values normalized
 582 to sum to 100 for interpretability, as well as confidence intervals (in brackets) for precision, recall,
 583 F1-score, and accuracy. As previously discussed, there exists a statistically significant gap between
 584 the model’s discriminative performance on the state prediction task and its actual success rate in
 585 generating correct completions. These findings suggest a promising direction for leveraging the
 586 model’s discriminative strength to improve its generative behaviour, which has been explored by Guo
 587 et al. [35].

Actual	Predicted		Metrics
	Positive	Negative	
Positive	25.40	3.56	Precision: 90.00%(86.02% – 93.78%)
Negative	2.82	68.22	Recall: 87.71%(83.46% – 91.47%)
Accuracy	93.62%(91.90% – 95.21%)		F1-score: 88.84%(84.72% – 92.61%)

Table 3: Confusion matrix of predicting next states by the model π^{sft} SFTed on both function call and state prediction data.

588 As discussed in Section 4.2, we track precision, recall, F1-score, and accuracy throughout the course
 589 of online RL training. The corresponding curves are shown in Figure 6. The red curves represent
 590 the model π^{sft} , which is initialized from SFT on $\mathbb{D}_{\text{fc}} \cup \mathbb{D}_{\text{sp}}$, while the blue curves correspond to the
 591 baseline model ϕ^{sft} , fine-tuned only on \mathbb{D}_{fc} . As shown in the figure, ϕ^{sft} , which lacks initial state
 592 prediction capability, consistently underperforms π^{sft} across all metrics throughout training. Even
 593 after 600 steps of RL training, ϕ^{rd} fails to match the performance of the SFT-only model π^{sft} , which
 594 indicates the necessary of the state prediction data. These results suggest that the benefits of \mathbb{D}_{sp}
 595 cannot be compensated for by relying solely on the downstream DyMo loss during RL training.

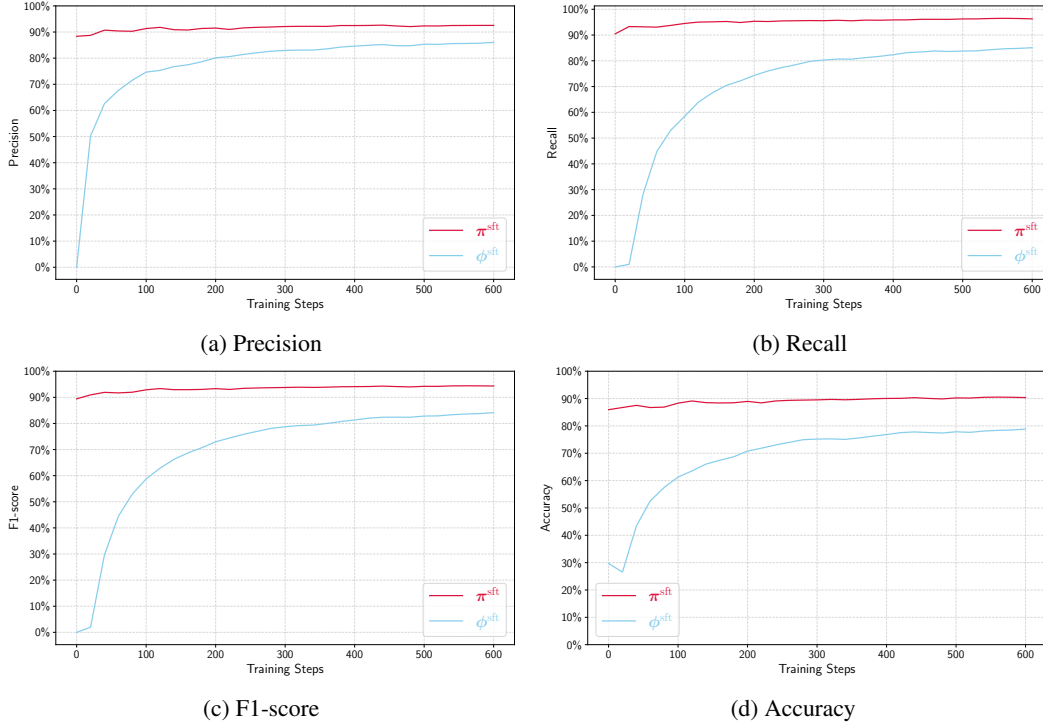


Figure 6: Performance metrics of results prediction over online RL training with DM loss (Equation 4).

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: Each key experimental finding is clearly articulated and corresponds directly to the claims made, ensuring a coherent and faithful representation of the work’s outcomes.

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: We discuss the limitations of our work and outline potential future directions to address them in Section 5.

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: [NA]

Justification: We do not include theoretical results in this work.

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 cross-referenced.
- 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: We provide comprehensive details of our experimental setup, ensuring that our main results can be reproduced. Additionally, we will release the simulator used in our experiments to further facilitate reproducibility.

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: Due to proprietary restrictions, we are unable to share the dataset used in our experiments. However, to facilitate reproducibility, we will release the simulator code employed in our study, along with comprehensive documentation and instructions. This will enable others to replicate our experimental setup and validate the main findings of our paper.

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: We provide comprehensive details of our training and testing procedures in Sections 3 and 4. This includes information on data splits, hyperparameter settings and selection strategies, optimizer types, and other relevant configurations. These details are provided to ensure that readers can fully understand and, if desired, replicate our experimental results.

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: In Section 4.2, we report confidence intervals for the binary classification experiments, which are computationally less intensive. These intervals were calculated using bootstrapping, capturing variability due to sampling randomness. However, for more computationally demanding experiments presented in other sections, we did not include error bars due to resource constraints. We acknowledge this limitation and plan to provide comprehensive statistical significance analyses for these experiments in future work.

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: [\[No\]](#)

Justification: Due to space constraints in the main body, we were unable to include detailed information about the computational resources used for each experiment. However, we acknowledge the importance of such details for reproducibility.

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: We have diligently ensured that our research aligns with the NeurIPS Code of Ethics. Our study does not involve human subjects or sensitive personal data, thereby avoiding concerns related to privacy, consent, or bias.

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: [No]

Justification: Our research focuses on enhancing the capabilities of LLMs through tool use—a topic that has been extensively explored in prior studies. Given this context, the societal impacts of our work are largely influenced by the specific tools integrated with the LLMs rather than our methodology itself. While our approach aims to improve the efficiency and versatility of LLMs, we acknowledge that the broader societal implications, both positive and negative, are contingent upon the applications of these tools. We encourage practitioners to consider these factors when deploying such systems.

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 research utilises proprietary datasets and models that, due to licensing restrictions, cannot be publicly released. Consequently, there is no risk of misuse associated with their release. The components we do share, such as simulator code, are built upon existing tools that have been extensively studied and are not associated with high-risk misuse scenarios. Therefore, we believe this question is not applicable to our work.

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 ensured that all third-party assets utilised in our research are properly credited. In the code repository accompanying our paper, we explicitly acknowledge the original repositories and their authors. Additionally, we have reviewed and adhered to the licenses and terms of use associated with these assets, ensuring compliance with their respective conditions.

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: [Yes]

Justification: We have released our customised simulator code under the same license as the original repository, following the original terms of use. To facilitate reproducibility and ease of use, we provide comprehensive documentation alongside the code, including detailed setup instructions, usage examples, and configuration guidelines. This documentation is included in the code repository and is designed to assist users in effectively utilising the simulator.

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: We do 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: We do 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.

- 936 • We recognize that the procedures for this may vary significantly between institutions
937 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
938 guidelines for their institution.
- 939 • For initial submissions, do not include any information that would break anonymity (if
940 applicable), such as the institution conducting the review.

941 16. Declaration of LLM usage

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

946 Answer: [NA]

947 Justification: In our research, we utilised LLMs solely for the purpose of refining the
948 language and grammar of the manuscript. These models were not employed in any aspect of
949 the research methodology, data analysis, or experimental design. NeurIPS

950 Guidelines:

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