# 🌮 Steering Vision-Language-Action Models as Anti-Exploration: A Test-Time Scaling Approach

**Siyuan Yang**[1,2,*], **Yang Zhang**[1,3,*,‡], **Haoran He**[4], **Ling Pan**[4], **Xiu Li**[3], **Chenjia Bai**[1,†], **Xuelong Li**[1,†]

[1]Institute of Artificial Intelligence, China Telecom, [2]University of Science and Technology of China, [3]Tsinghua University, [4]The Hong Kong University of Science and Technology
[*]Equal Contributions   [†]Corresponding Authors   [‡]Project Leader

Vision-Language-Action (VLA) models, trained via flow-matching or diffusion objectives, excel at learning complex behaviors from large-scale, multi-modal datasets (e.g., human teleoperation, scripted policies). However, since VLAs incorporate diverse data modes in the pre-training stage, and the finetuning dataset often contains demonstration data collected in a kinematically suboptimal or undesirable way, it exists redundant action modes that are irrelevant to the success action modes of the downstream task. Specifically, we observe a critical inference-time fragility among various sampled noises after supervised finetuning of pre-trained VLAs. In this paper, we attribute this instability to the distribution shift between the VLA policy and the policy induced by stable success modes of the downstream task dataset. Thus, we propose **TACO**, a test-time-scaling (TTS) framework that applies a lightweight pseudo-count estimator as a high-fidelity verifier of action chunks. The VLA models integrated with TACO can execute the actions with maximum pseudo-count from all sampled action chunks, thereby preventing distribution shifts while preserving the generalization ability of VLAs since the constraint is applied only during inference. Our method resembles the classical anti-exploration principle in offline reinforcement learning (RL), and being gradient-free, it incurs significant computational benefits compared to RL update, especially for flow or diffusion-based VLAs which are difficult to perform RL update due to denoising process. Extensive experiments across four simulation benchmarks (RoboTwin2.0, Robotwin, LIBERO, SimplerEnv) and a dual-arm platform demonstrate that our method significantly improves the inference stability and success rates in downstream-task adaptations.

🛡️TeleAI **TeleAI**

## 1   Introduction

Vision-Language-Action (VLA) models (Brohan et al., 2022; Zitkovich et al., 2023; Kim et al., 2024; Black et al., 2024) have demonstrated remarkable performance and strong generalization capabilities in robotic manipulation tasks, benefiting from large-scale, multimodal datasets collected via human teleoperation or scripted policies. During pretraining, VLA models that employ flow-matching (Liu et al., 2023b; Lipman et al., 2023, 2024) or diffusion (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2021) -based objectives can effectively learn multimodal behavioral policies that span diverse manipulation modes, scenes, and tasks (Jang et al., 2021; Collaboration., 2023; Walke et al., 2023; Khazatsky et al., 2024). Subsequently, the pretrained VLA model is adapted to downstream tasks through supervised fine-tuning (SFT) on small-scale, task-specific datasets. Leveraging a pretrained VLA for downstream adaptation typically yields substantially better generalization ability than training a policy from scratch, as the model benefits from rich pretrained knowledge—including visual perception, language understanding, embodiment priors, and awareness of object types and spatial configurations (Driess et al., 2023; Beyer et al., 2024; Wang et al., 2024).

Although VLA models exhibit strong average performance, we observe a critical fragility at inference time in representative flow-matching (Black et al., 2024) and diffusion-based VLA (Liu et al., 2025) architectures—even

after SFT on task-specific datasets (Mu et al., 2025). As illustrated in Figure 1, we evaluate the same SFT-adapted VLA model using two fixed noise vectors during inference across different tasks and find that the resulting success rates vary drastically—from 0% to 80%—solely due to noise variation.

We attribute this instability to the persistence of redundant action modes that are irrelevant to the success action modes of the downstream task after SFT. Specifically, (i) during pretraining, VLA models absorb a broad spectrum of action modes from diverse data sources, making it difficult to rapidly narrow their output distribution to the narrow set of successful behaviors required by a specific downstream task. Consequently, the policy distribution after SFT still retains extraneous modes unrelated to task success. (ii) SFT datasets themselves may exhibit multimodality, as they are often collected from multiple human teleoperators, scripted planners, or varying execution styles—some of which encode suboptimal or undesirable strategies. These redundant modes induce a significant distributional shift between the VLA policy and the ideal policy corresponding to the stable success mode of the downstream task. This shift becomes evident when sampling different noise vectors in flow-matching or diffusion-based frameworks, as the stochasticity of the sampled initial noise directly influences the resulting actions.
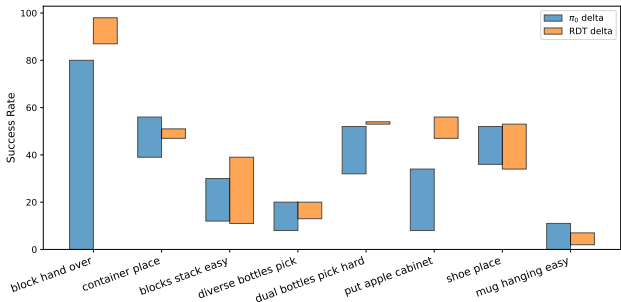


**Figure 1 Illustrative Example:** We fix the initial noise for the denoising process of $\pi_0$ to $noise_1$ and evaluate the success rate. Then we fix the initial noise for the denoising process of $\pi_0$ to another noise value $noise_2$ and evaluate the success rate under the same scenarios. We plot the two resulting success rates using a floating bar chart. The same procedure is applied to RDT.

In this paper, we address the distribution shift problem from an anti-exploration perspective and propose **T**est-time **A**nti-exploration via pseudo-**CO**unts (**TACO**). The anti-exploration principle (Rezaeifar et al., 2022) originates from offline reinforcement learning (RL) (Levine et al., 2020b), where the goal is to prevent the policy from visiting states or actions that lie outside the support of the dataset. Analogously, during VLA inference, we aim to constrain the generated actions to lie within the support of the successful modes present in the SFT dataset—avoiding exploration of redundant or irrelevant action modes retained from pretraining or imperfect fine-tuning data. In TACO, we realize anti-exploration through test-time scaling (TTS) rather than policy optimization, as flow-matching or diffusion-based VLA models involve complex sampling dynamics that are not easily amenable to standard RL-style updates. Instead, TTS adjusts the action sampling process without modifying the VLA's parameters, thus requiring no gradient computation. We employ a classical Coin-Flipping Network (CFN) (Lobel et al., 2023) enhanced with an internal representation mechanism to estimate pseudo-counts for every observation-instruction-action-chunk pair in the SFT dataset. Intuitively, an action chunk with a higher pseudo-count is more consistent with the frequently observed (i.e., successful) behaviors in the SFT data and is thus preferred. During inference, TACO uses the CFN-derived pseudo-counts as a verifier to select the most reliable action chunk among multiple action chunk candidates, ultimately executing the one with the highest pseudo-count. To mitigate the computational overhead of repeated forward passes at test-time, we further introduce a shared observation key-value cache, which dramatically reduces latency by reusing visual-language representations across samples.

Our contributions can be summarized as follows:

- We propose TACO, a test-time scaling framework for VLAs which retains the strong generalization capabilities of pretrained VLAs while effectively constraining outputs to the success modes of specific downstream tasks.

- We introduce an efficient internal representation mechanism for pseudo-count estimation in CFN, enabling accurate measurement of distributional shift with minimal computational overhead.

- We demonstrate that TACO significantly boosts the success rates of diverse VLA models across both simulation and real-world tasks without prolonged training time and can operate with low latency.

## 2 Related Work

**Vision-Language-Action Models.** Vision-Language-Action (VLA) models have emerged as a dominant paradigm for general-purpose robotic manipulation, integrating the reasoning capabilities of pre-trained VLMs with robotic control. While early works like RT-1 and Octo (Brohan et al., 2022; Team et al., 2024) trained policies from scratch, later approaches such as RT-2 and OpenVLA (Zitkovich et al., 2023; Kim et al., 2024) successfully leveraged pre-trained backbones (e.g., PaLI-X, Prismatic) (Chen et al., 2023; Anschütz and Le Bras, 2023) via action discretization. To better capture the multimodality and precision of action distribution in large-scale cross-embodiment action-labeled datasets (Collaboration., 2023; Wu et al., 2024; Bu et al., 2025), recent approaches have shifted towards generative policies. Methods like RDT-1B, DexVLA, DexGraspVLA, Dita (Liu et al., 2025; Wen et al., 2025; Zhong et al., 2025; Hou et al., 2025) introduced diffusion (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2021) architectures, and the latest generation of models—including $\pi_0$, GR00T N1, and SmolVLA (Black et al., 2024; Bjorck et al., 2025; Shukor et al., 2025)—leverage Mixture-of-Transformers architecture and employ flow matching (Liu et al., 2023b; Lipman et al., 2023, 2024) or diffusion experts to synthesize continuous action sequences. Addressing the distribution shift in these score-based generative frameworks, recent work has also explored latent alignment. For instance, ATE (Zhang et al., 2025c) introduces a unified latent space to align disparate action distributions and employs classifier guidance during training time, to steer denoising-based VLAs for efficient adaptation. Despite this rapid architectural progress, a critical and often overlooked issue is the inference-time instability of these generative models. Our work is the first to systematically diagnose this instability as an out-of-support problem and propose a test-time correction mechanism based on anti-exploration.

**Anti-exploration & Offline RL.** In online RL (Sutton et al., 1999), the agent actively interacts with the environment to explore the state–action space (Yang et al., 2021). In contrast, offline RL learns a policy exclusively from a fixed, previously collected dataset, without further environmental interaction (Levine et al., 2020a; Rashidinejad et al., 2021). This setting poses significant challenges—even for off-policy methods—due to distributional shift between the behavior policy that generated the data and the target policy being learned (Wu et al., 2019). Specifically, the offline dataset reflects the visitation distribution of the behavior policy, which may differ substantially from that of the learned policy (Agarwal et al., 2019). This discrepancy can cause severe extrapolation errors when estimating values for out-of-distribution (OOD) state–action pairs (Kumar et al., 2020; Bai et al., 2022). A similar distribution shift arises during the SFT stage of VLAs, where the VLA model encode a broad and redundant action modes, whereas successful behaviors in downstream tasks typically occupy a much narrower region. Consequently, when sampling actions in flow or diffusion-based process, the policy may generate suboptimal behaviors (Lipman et al., 2023; Ho et al., 2020). In TACO, we employ an anti-exploration strategy (Rezaeifar et al., 2022; Nikulin et al., 2023) that imposes policy constraints to restrict the policy to stay close to the support of the success data distribution. Unlike conventional offline RL approaches that typically rely on uncertainty-aware or conservative penalties (Fujimoto and Gu, 2021; Kidambi et al., 2020; An et al., 2021), our method adopts a simple count-based mechanism (Ostrovski et al., 2017; Bellemare et al., 2016; Lobel et al., 2023), which is well-suited to the structure of VLA models and can operate without of value estimation (Kostrikov et al., 2022). Notably, count-based regularization has precedent in both classical RL and recent RLHF algorithms.

**Test Time Scaling in VLA.** Performance improvements at test time by extending inference have been widely validated in LLMs (Ehrlich et al., 2025; Chen et al., 2024; Brown et al., 2024). In VLA research, many works achieve performance gains through internal 'thinking' mechanisms, but most of them require datasets with annotated reasoning traces (Zhao et al., 2025). In contrast, test-time scaling introduces an additional scoring module to select better candidates, enabling a generate-then-justify mechanism to further boost the performance of the policy without modifying the network weights, as demonstrated by methods that utilized the advantage function (Zhang et al., 2025d) and state-action value function (Nakamoto et al., 2024), and ohther variants like Verifier-free (Jang et al., 2025) that considered the likelihood. Recent work (Du and Song, 2025) also explored leveraging the distance to the goal within the representation of a learned world model (Hafner et al., 2019, 2020; Zhang et al., 2025a,b) as a metric to steer the base policy. Repeated sampling combined with a verifier has also been shown effective in several works (Song et al., 2021; Kwok et al., 2025); however, these approaches typically involve complex RL training or rely on verifiers with a large number of parameters. To the best of our knowledge, no prior work has leveraged the internal features of a VLA to

build a lightweight and efficient verifier that fully exploits the model's understanding capabilities, without modifying the backbone network.

## 3  Preliminaries

**Problem Statement.** We consider a language-conditioned robotic manipulation task with vision input under the imitation learning setting. We assume access to a pre-collected dataset $\mathcal{D}_{\mathrm{sft}} = \{\tau_1, \tau_2, \ldots, \tau_n\}$ comprising $n$ demonstrations from one or more target downstream tasks. Each demonstration $\tau = (\mathbf{o}_{1:T}, l, \mathbf{a}_{1:T})$ contains a natural language instruction $l$ describing the target task, a sequence of observations $\mathbf{o}_{1:T}$, where each observation $\mathbf{o}_t = (\mathbf{I}_t^1, \ldots, \mathbf{I}_t^n, \mathbf{q}_t)$ consists of multiple RGB images $(\mathbf{I}_t^1, \ldots, \mathbf{I}_t^n)$ and the robot's proprioceptive state $\mathbf{q}_t \in \mathbb{R}^m$, and a sequence of actions $\mathbf{a}_{1:T}$ where $\mathbf{a}_t \in \mathbb{R}^n$. In practice, expert demonstrations in $\mathcal{D}_{\mathrm{sft}}$ are scarce due to the high cost of obtaining them. Thus, they typically cover a limited fraction of the whole observation and action spaces, which in turn makes the distribution shift issue potentially more pronounced.

**Coin Flipping Network.** Coin Flipping Network (CFN) (Lobel et al., 2023) is a neural estimator which uses the sampling distribution of Rademacher trials (or *coin flips*) made every time a state is encountered to derive the state visitation counts. Given a state $s$, each occurrence of $s$ is paired with a random binary vector $\mathbf{c}_i \sim \{-1, 1\}^d$, forming a dataset $\mathcal{D}_{\mathrm{cfn}} = \{(s_i, \mathbf{c}_i)\}$. The CFN $f_\phi$ parameterized by $\phi$ is learned by solving the following simple regression problem:

$$\min_\phi \mathbb{E}_{(s_i, \mathbf{c}_i) \sim \mathcal{D}_{\mathrm{cfn}}} \left[ \| f_\phi(s_i) - \mathbf{c}_i \|^2 \right]. \tag{1}$$

For a state $s$ appearing $m$ times in $\mathcal{D}_{\mathrm{cfn}}$, the optimal solution satisfies $f_\phi^*(s) = \frac{1}{m} \sum_{i=1}^m \mathbf{c}_i$. If each element $c_{ij}$ in $\mathbf{c}_i$ follows a fair coin flipping distribution,

$$\frac{1}{d} \| f_\phi(s) \|^2 = \frac{1}{d} \sum_{j=1}^d \mathbb{E}\left[ (\sum_{i=1}^m \frac{c_{ij}}{m}) \right] = \frac{1}{m}. \tag{2}$$

Thus, we can simply approximate the inverse of the state visitation count given by $\| f_\phi(s) \|^2 / d \approx 1/N(s)$. In this work, the CFN would be utilized for *anti-exploration*.

## 4  Method

We now present TACO, a framework that treats the pseudo count estimator as an off-the-shelf verifier to scale test-time compute for VLAs, realizing *Anti-Exploration* principle. Note that TACO can be directly integrated into VLAs that either use discrete token auto-regression (e.g., OpenVLA (Kim et al., 2024)) or a score-based generative formulation (e.g., RDT (Liu et al., 2025), $\pi_0$ (Black et al., 2024), $\pi_{0.5}$ (Black et al., 2025)) to model the action distribution. We first elucidate how we draw inspiration from *Anti-Exploration* in offline RL to mitigate the distribution shift between the fine-tuned VLA and the desired success mode in the provided dataset of the downstream task (§4.1). Then, we describe in detail how we operationalize the principle of *Anti-Exploration* (§4.2) and incorporate it to drive VLAs from being out-of-support during inference time (§4.3). An overview of TACO is shown in Figure 2.

### 4.1  Inference Instability as an OOD Problem

We begin by formally identifying the source of inference-time instability (§1). VLAs trained via SFT (Black et al., 2024; Liu et al., 2025) are compelled to approximate the entire dataset distribution $p_\mathcal{D}(a|s)$. This dataset is often an imperfect mixture $p_\mathcal{D} = w_1 \pi^* + \sum_{k>1} w_k p_k$, where $\pi^*$ is the desired success mode and $p_{k>1}$ represent sub-optimal or unexpected modes. The resulting policy, $\pi_\theta \approx p_\mathcal{D}$, is thus inherently multimodal. Instability arises when $\pi_\theta$ generates an action $\mathbf{a}_j \sim p_{j>1}(\mathbf{a}|\mathbf{o}, l)$ from an undesired mode. However, $\pi_\theta$, trained via offline imitation learning (IL), is mode-agnostic and cannot distinguish successful from unsuccessful modes.

This paradigm mirrors the core challenge in Offline RL (Levine et al., 2020b): mitigating extrapolation error (Fujimoto et al., 2019; Kumar et al., 2019, 2020; Kostrikov et al., 2022). Generating an action from an undesired mode $p_{j>1}$ is analogous to an Offline RL agent sampling an action that is out-of-support relative to the target success mode $\pi^*$. A principled solution in Offline RL is *anti-exploration* (Rezaeifar et al., 2022),
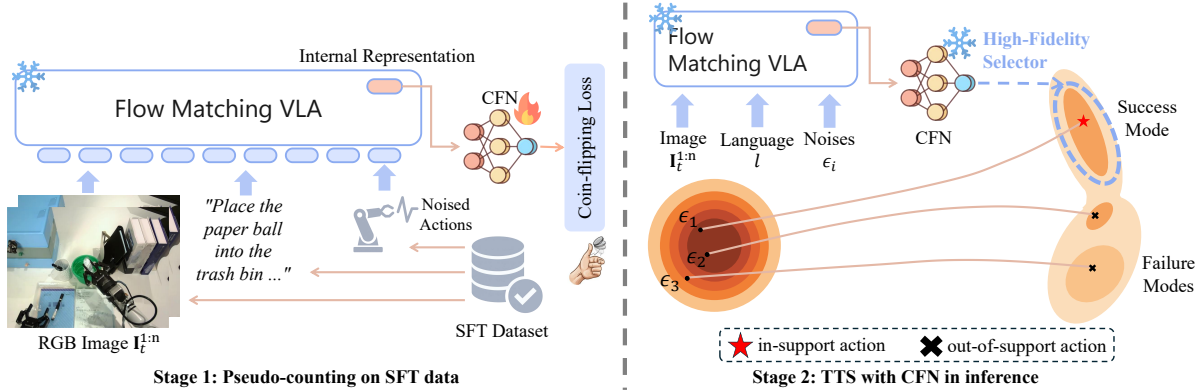
**Figure 2 Overview of TACO.** In the training stage (Stage 1), we sample data from the SFT dataset, add a certain amount of noise to the expert actions, and feed them into the VLA to denoise the actions while extracting internal representations. These representations are then used to train the CFN. During inference (Stage 2), the VLA generates multiple candidate actions along with their corresponding internal representations, and the CFN serves as a selector to select the action with the highest count for execution.

which optimizes a penalized value function to constrain the policy to stay within the support of the dataset. This provides our starting objective:

$$\mathbf{a}^* = \arg\max_{\mathbf{a}}[Q(s, a) - b(s, a)] \tag{3}$$

where $Q(s, a)$ is the action-value function and $b(s, a)$ is an anti-exploration penalty (or "bonus") that assigns high cost to out-of-support state-action pairs.

To relate this objective to our test-time inference framework, we first introduce a ideally simplified setting: we model the task as a contextual bandit where executing an action chunk $\mathbf{a}_{1:H}$ given $(\mathbf{o}, l)$ yields a binary success ($r = 1$) or failure ($r = 0$) reward. Under this formulation, the Q-function becomes equivalent to the success probability: $Q(\mathbf{o}, l, \mathbf{a}_{1:H}) = P(r = 1|\mathbf{o}, l, \mathbf{a}_{1:H})$. The anti-exploration objective thus translates to,

$$\mathbf{a}^*_{1:H} = \arg\max_{\mathbf{a}_{1:H}}[P(r = 1|\mathbf{o}, l, \mathbf{a}_{1:H}) - b(\mathbf{o}, l, \mathbf{a}_{1:H})]. \tag{4}$$

To solve Eq. (4), we introduce our core assumption: the downstream dataset $\mathcal{D}_{\mathrm{sft}}$ for VLA fine-tuning is demonstrative. That is, we assume the density of demonstrations correlates with task success; high-density modes correspond to an ideal policy $\pi^*$ and low-density modes correspond to $p_k$. Formally,

$$P(r = 1|\mathbf{o}, l, \mathbf{a}_{1:H}) \propto p_{\mathcal{D}_{\mathrm{sft}}}(\mathbf{o}, l, \mathbf{a}_{1:H}). \tag{5}$$

Let the anti-exploration penalty be a count-based bonus $b(\mathbf{o}, l, \mathbf{a}_{1:H}) = 1/N_{\mathcal{D}_{\mathrm{sft}}}(\mathbf{o}, l, \mathbf{a}_{1:H})$ where $N_{\mathcal{D}_{\mathrm{sft}}}(\mathbf{o}, l, \mathbf{a}_{1:H})$ is the corresponding visit count in the dataset $\mathcal{D}_{\mathrm{sft}}$.

This instantiation formally simplifies the principled Offline RL objective in Eq. (4) to a equivalent objective: finding the action chunk with the **maximum visitation count**, i.e., density:

$$\mathbf{a}^*_{1:H} = \arg\max_{\mathbf{a}_{1:H}}[N_{\mathcal{D}_{\mathrm{sft}}}(\mathbf{o}, l, \mathbf{a}_{1:H})]. \tag{6}$$

This result provides a strong theoretical justification for our method. Our test-time scaling framework, which seeks the highest-density (i.e., in-support) action, is not merely a heuristic. It constitutes a principled and theoretically equivalent realization of the anti-exploration objective standard in Offline RL.

## 4.2 Coupled Pseudo-Count Estimation for VLAs

To realize our anti-exploration principle (Eq. (6)), we must estimate the visit count $N_{\mathcal{D}_{\mathrm{sft}}}(\mathbf{o}_t, l, \mathbf{a})$[1]. We can replace it with a pseudo-count $\hat{N}_{\mathcal{D}_{\mathrm{sft}}}$ by training a pseudo-count estimator $f_\phi$ on $\mathcal{D}_{\mathrm{sft}}$. Motivated by recent

---

[1]In the following, we omit the subscript of action chunk for clarity.

successes (Lobel et al., 2023; Bai et al., 2025) on applying pseudo-counts to LLM preference optimization, we instantiate this estimator as a Coin Flipping Network (CFN).

**The Coupled Estimator.** A critical design choice is the input representation $z = \text{Enc}(\mathbf{o}, l, \mathbf{a})$. Instead of training a separate encoder, we hypothesize that the VLA $\pi_\theta$ itself provides the richest representations. We instantiate the CFN $f_\phi$ as a simple, lightweight MLP head that takes the VLA's internal representation $h_\theta$ as input (see Figure 2). This coupled estimator design is highly efficient, leveraging the VLA's computation, and benefits from its extensive pre-training.

**High-Fidelity Feature Search for Denoising-based VLAs.** This coupled design introduces a critical challenge for diffusion-based or flow-based VLAs (Black et al., 2024; Liu et al., 2025; Black et al., 2025). These models are trained exclusively on noised actions $\{\mathbf{a}_\sigma\}$ and thus never see clean data actions $\{\mathbf{a}\}$ during training. Directly extracting $h_\theta$ by feeding $\mathbf{a} \in \mathcal{D}_{\text{sft}}$ may not lie on VLA's feature space, yielding uninformative representations.

Our goal is to find an in-distribution feature $h_\theta$ that best represents the clean action $a$. We propose a **High-Fidelity Feature Search** procedure. For each data point $(\mathbf{o}, l, \mathbf{a}) \in \mathcal{D}_{\text{sft}}$, we query the VLA $N$ times with different noise levels $\{\sigma_i\}_{i=1}^N$:

$$\{(\mathbf{a}_{\text{pre}}^{(i)}, h_\theta^{(i)})\}_{i=1}^N = \{\text{VLA}(\mathbf{o}, l, \mathbf{a}_{\sigma_i})\}_{i=1}^N \tag{7}$$

where $\mathbf{a}_{\sigma_i}$ is the corrupted version of the original action, $\mathbf{a}_{\text{pre}}^{(i)}$ is the VLA's resulting clean prediction, and $h_\theta^{(i)}$ is the corresponding internal representation. We then select the feature $h_\theta^{(i^\star)}$ whose corresponding action prediction $\mathbf{a}_{\text{pre}}^{(i^\star)}$ is closest to the ground-truth $\mathbf{a}$:

$$i^\star = \arg\min_{i \in \{1, \ldots, N\}} \|\mathbf{a}_{\text{pre}}^{(i)} - \mathbf{a}\|_2. \tag{8}$$

The resulting feature $h_\theta^{(i^\star)}$ is now both *in-distribution* for the VLA (as it was generated from a noised input $\mathbf{a}_{\sigma_{i^\star}}$) and *high-fidelity* (as it is confirmed to represent $\mathbf{a}$).

**Estimator Training.** This search procedure yields a high-fidelity feature set $\mathcal{D}_h = \{h_\theta^{(i^\star)}\}$ representing our original dataset $\mathcal{D}_{\text{sft}}$. We can now train the CFN head $f_\phi$ on top of any VLA's last hidden states. Following (Lobel et al., 2023), the CFN $f_\phi$ is learned via optimizing Eq. (1). After training the CFN, the pseudo-count of a feature $h_\theta$ is naturally given by

$$\hat{N}_{\mathcal{D}_{\text{sft}}}(\mathbf{o}, l, \hat{\mathbf{a}}) = \hat{N}_{\mathcal{D}_{\text{sft}}}(h_\theta) \propto \frac{1}{\|f_\phi(h_\theta)\|^2}. \tag{9}$$

where $h_\theta$ is the feature extracted from a candidate action $\hat{\mathbf{a}}$ during inference.

As empirically validated in Figure 3, the pseudo-counts derived from our coupled estimator exhibit a strong correlation between the L2 distance to the ground-truth action and the pseudo-count estimated by our CFN. This confirms that our coupled CFN can effectively function as the verifier needed to penalize out-of-support actions.

## 4.3 Test-Time Scaling as Anti-Exploration

We now instantiate our anti-exploration framework through a two-stage **generate-then-verify** procedure during inference. First, in the *generation stage*, given an obsevation $\mathbf{o}_t$ and a language instruction $l$, we leverage the VLA $\pi_\theta$ which is already fine-tuned in the downstream dataset $\mathcal{D}_{\text{sft}}$, as a candidate generator. For diffusion or flow-based models, this is achieved by sampling $M$ distinct initial noise vectors $\{\epsilon_i\}_{i=1}^M$ and performing the full denoising process in a batch-parallel way, yielding a batch of $M$ candidate action chunks $\{\hat{\mathbf{a}}_{t:t+H}^{(i)}\}_{i=1}^M$. Crucially, we simultaneously extract their corresponding internal representations $\{h_\theta^{(i)}\}_{i=1}^M$ during this generation pass.

Second, in the *verification stage*, we deploy our trained CFN $f_\phi$ (in §4.2) as the verifier. The CFN scores each candidate $\hat{\mathbf{a}}_{t:t+H}^{(i)}$ based on its feature $h_\theta^{(i)}$, effectively estimating its pseudo-count $\hat{N}_{\mathcal{D}_{\text{sft}}}(\mathbf{o}_t, l, \hat{\mathbf{a}}_{t:t+H}^{(i)}) = 1/\|f_\phi(h_\theta^{(i)})\|^2$. Finally, we select the single action $\hat{\mathbf{a}}_{t:t+H}^*$ that maximizes this pseudo-count,

$$\hat{\mathbf{a}}_{t:t+H}^* = \hat{\mathbf{a}}_{t:t+H}^{(i^*)}, \text{ where } i^* = \arg\max_{i \in \{1, \ldots, M\}} \frac{1}{\|f_\phi(h_\theta^{(i)})\|^2}.$$

This procedure is a direct realization of our principled objective in Eq. (6). Instead of taking a random sample from the VLA's potentially unstable multimodal distribution, we deterministically select the action that is most **in-support** (i.e., highest-count) w.r.t. the successful modes of $\mathcal{D}_{\text{sft}}$.

**Enabling Efficient Inference via KV Cache Optimization.** Naively sampling $M$ candidates introduces a prohibitive $\mathcal{O}(M)$ overhead, rendering our method impractical. We solve this by observing that the expensive VLA computations (e.g., the transformer backbone) depend *only* on the shared context tokens $(\mathbf{o}, l)$. Therefore, we compute the Key (K) and Value (V) caches for this shared context just once and re-use them across all $M$ parallel action generation (e.g., denoising) processes. This optimization is a key of our method, making the marginal cost of additional candidates minimal. As shown in Figure 4 (left), when sampling with $M = 32$, our method reduces inference time by **73.2%** compared to naive multi-batch parallel inference, making our high-performance anti-exploration sampling practical. We provide an algorithmic description on how TACO does test-time scaling in §A.

# 5    Experiments

Our experiments aim to verify whether the proposed Test-time Anti-exploration via pseudo-COunts can effectively improve the accuracy of action prediction, thereby enhancing task success rates. We also seek to understand which design choices contribute to these improvements, and whether the additional computational cost is acceptable. Specifically, our experiments are designed to answer the following three research questions:

**Q1:** Can our framework improve performance across diverse environments, various VLA policies, and a wide range of tasks?

**Q2:** What is the additional time cost introduced by our test-time scaling method? Can it run efficiently on real robots while maintaining performance gains?

**Q3:** How do design choices—such as the use of internal features and the CFN pseudo-count network—affect overall performance?

To this end, we evaluate our framework on four simulation benchmarks with 64 tasks and 5 real-world tasks, covering three types of VLA policies.

## 5.1    Simulation Experiments

### 5.1.1   Setup and Baselines

**Benchmarks.** We evaluate our test-time scaling framework on the *Simpler* (Li et al., 2024b), *Libero* (Liu et al., 2023a), *Robotwin1.0* (Mu et al., 2025), and *Robotwin2.0* (Chen et al., 2025) benchmarks. *Simpler* is a real-to-sim benchmark built upon the SAPIEN simulator and the ManiSkill2 benchmark. It provides simulated task environments for both the WindowX and Google Robot platforms; in this work, we primarily utilize the tasks designed for WindowX. *Libero* is a benchmark for lifelong learning in decision making (LLDM), consisting of four task suites. We focus on the most challenging suite, as previous works have already achieved near-perfect success rates on the others. *Robotwin1.0* and *Robotwin2.0* mainly focus on dual-arm gripper manipulation, featuring a large variety of assets and task types, and providing scripted policies for automatic data collection.

**Baselines.** We primarily evaluate our framework on flow-matching-based VLA policies, including $\pi_0$ (Black et al., 2024) and $\pi_{0.5}$ (Black et al., 2025). To demonstrate the generality of our framework, we also apply test-time scaling to the autoregressive-based *OpenVLA* (Kim et al., 2024) framework and compare it with *RoboMonkey* (Kwok et al., 2025), which is likewise a test-time scaling framework. *RoboMonkey* builds upon *LLaVA-7B* (Touvron et al., 2023) and employs preference learning to train an action verifier. For a more comprehensive comparison, we also report the success rates of RT-1-X (O'Neill et al., 2024), Octo (Team et al., 2024), RoboVLM (Li et al., 2024a), SpatialVLA (Qu et al., 2025), and RDT (Liu et al., 2025) on selected benchmarks.

**Table 1** Evaluation of success rate(%) on the RoboTwin 1.0.

| | Block Handover | Bottles Adjust | Container Place | Diverse Bottles Pick | Dual Bottles Pick Easy |
|---|---|---|---|---|---|
| $\pi_0$ | 41.0 | 31.0 | 25.0 | 21.0 | 60.0 |
| $\pi_0$ **+ TACO (Ours)** | **62.0**(↑ 21) | **40.0**(↑ 9) | **40.0**(↑ 15) | **27.0**(↑ 6) | **70.0**(↑ 10) |

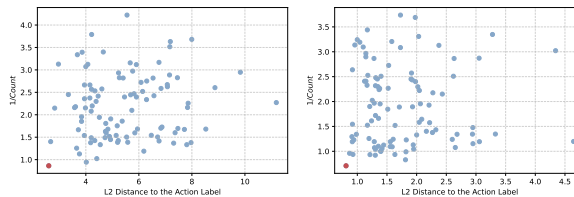| | Dual Bottles Pick Hard | Pick Apple Messy | Shoe Place | Mug Hanging Easy | Average |
|---|---|---|---|---|---|
| $\pi_0$ | 48.0 | 15.0 | 42.0 | 7.0 | 32.2 |
| $\pi_0$ **+ TACO (Ours)** | **52.0**(↑ 4) | **19.0**(↑ 4) | **50.0**(↑ 8) | **12.0**(↑ 5) | **41.3**(↑ 9.1) |

**Table 2** Success rate (%) evaluated on Simpler-WindowX. Baseline results are taken from *Simpler*, *RoboVLM*, and *SpatialVLA*. Our method achieves an average improvement of 7.5% over $\pi_0$.

| | RT-1-X | Octo | RoboVLM | SpatialVLA | $\pi_0$ | $\pi_0$ **+ TACO (Ours)** |
|---|---|---|---|---|---|---|
| Spoon on Towel | 0.0 | 12.5 | 29.2 | 16.7 | 36.0 | **52.0**(↑ 16) |
| Carrot on Plate | 4.2 | 8.3 | 25.0 | 25.0 | 42.0 | **52.0**(↑ 10) |
| Stack Cubes | 0.0 | 0.0 | 12.5 | 29.2 | **34.0** | 30.0(↓ 4) |
| Eggplant in Basket | 0.0 | 43.1 | 58.3 | **100.0** | 80.0 | **88.0**(↑ 8) |
| Average | 1.1 | 16.0 | 31.3 | 42.7 | 48.0 | **55.5**(↑ 7.5) |

### 5.1.2 Results

To answer **Q1:** *Can our framework improve performance across diverse environments, various VLA policies, and a wide range of tasks?* We evaluate our test-time scaling framework on four benchmarks, all of which achieve the highest average success rate, as shown in Table 1 2 3. Specifically, our method improves performance by 9.1%, 7.5%, 4.7%, respectively. In the Libero benchmark, as shown in Table 4, since our base policy $\pi_{05}$ has been fine-tuned and achieved near-perfect success rates on three of the suites, we only evaluate our performance on the Libero-long tasks. Even though the base policy already reaches an average success rate of 94.8%, TACO further improves it by 1.8%. In particular, the success rate on the *Moka Pots on Stove* task increases from 68% to 86%. Moreover, we also evaluate our method when applied to OpenVLA in Libero. Despite only introducing a lightweight MLP pseudo-counter, it still helps OpenVLA achieve a 6.0% improvement in average success rate. For comparison, Robomonkey, trained with a 7B VLM, only achieves a 6.7% improvement.

To further investigate the underlying reason for the effectiveness of our approach, we conducted an additional analysis. For the same test set observations, we generated 100 random noise inputs and examined the relationship between the norm of the CFN (the reciprocal of the Count) output and the $L_2$ distance between the predicted action and the ground-truth action. As shown in Figure 3, the results indicate that selecting the action corresponding to the smallest CFN output norm (i.e., the largest Count) almost always yields the action closest to the ground truth, while effectively filtering out overly aggressive actions.



**Figure 3** We performed 100 action samples for a single observation from the test set in each of the *stamp seal*(left), *move can pot*(right) tasks, and scored them using the pseudo-counter. The red points indicate the actions with the highest Count values, which generally correspond to the actions with the smallest $L_2$ distance to the action labels.

## 5.2 Real-World Experiments

### 5.2.1 Setup

We build our real-world experimental setup using a RealMan75 dual-arm robot in an office-like environment. The setup includes five tasks involving various common objects such as books, pens, a cabinet, a charger, a laptop, and a paper ball. The five tasks are as follows: *Receive Book*, where the robot receives a book handed over by a human and places it back on the bookshelf; *Storage Charger*, where it picks up the charger head, places it inside the cabinet, and closes the cabinet door; *Paper and Pen*, where it picks up the paper ball and the pen and places them into the trash bin and pen holder, respectively; *Laptop*, where it closes the laptop lid,

**Table 3** Evaluation on the Robotwin2.0 benchmark. Each task is tested across 100 randomly generated scenes using 100 different seeds. For test-time scaling, the number of candidate actions is set to 50.

| | Adjust Bottle | Beat Block Hammer | Blocks Ranking Size | Dump Bin Bigbin | Grab Roller | Handover Block | Handover Mic |
|---|---|---|---|---|---|---|---|
| RDT | 81.0 | 77.0 | 0.0 | 64.0 | 74.0 | **45.0** | **90.0** |
| $\pi_{0.5}$ | 89.0 | 69.0 | 36.0 | 86.0 | **100.0** | 24.0 | 58.0 |
| $\pi_{0.5}$ + TACO (Ours) | **93.0** (↑ 4) | **79.0** (↑ 10) | **40.0** (↑ 4) | **87.0** (↑ 1) | 98.0 (↓ 2) | 36.0 (↑ 12) | 63.0 (↑ 5) |

| | Lift Pot | Move Can Pot | Move Pillbottle Pad | Move Playingcard Away | Move Stapler Pad | Open Laptop | Open Microwave |
|---|---|---|---|---|---|---|---|
| RDT | **72.0** | 25.0 | 8.0 | 43.0 | 2.0 | 59.0 | **37.0** |
| $\pi_{0.5}$ | 62.0 | 42.0 | **55.0** | 85.0 | 13.0 | **67.0** | 20.0 |
| $\pi_{0.5}$ + TACO (Ours) | 66.0 (↑ 4) | **57.0** (↑ 15) | 54.0 (↓ 1) | **88.0** (↑ 3) | **18.0** (↑ 5) | **67.0** (0) | 21.0 (↑ 1) |

| | Pick Diverse Bottles | Pick Dual Bottles | Place A2B Left | Place A2B Right | Place Bread Basket | Place Bread Skillet | Place Burger Fries |
|---|---|---|---|---|---|---|---|
| RDT | 2.0 | 42.0 | 3.0 | 1.0 | 10.0 | 5.0 | 50.0 |
| $\pi_{0.5}$ | 52.0 | 72.0 | 51.0 | 39.0 | 62.0 | **62.0** | 89.0 |
| $\pi_{0.5}$ + TACO (Ours) | **59.0**(↑ 7) | **73.0**(↑ 1) | **56.0**(↑ 5) | **42.0**(↑ 3) | **65.0**(↑ 3) | 61.0 (↓ 4) | **92.0**(↑ 3) |

| | Place Cans Plasticbox | Place Container Plate | Place Dual Shoes | Place Fan | Place Mouse Pad | Place Object Basket | Place Object Scale |
|---|---|---|---|---|---|---|---|
| RDT | 6.0 | 78.0 | 4.0 | 12.0 | 1.0 | 33.0 | 1.0 |
| $\pi_{0.5}$ | 68.0 | 85.0 | 23.0 | **34.0** | 16.0 | 69.0 | 45.0 |
| $\pi_{0.5}$ + TACO (Ours) | **72.0** (↑ 4) | **87.0** (↑ 2) | **28.0** (↑ 5) | 32.0 (↓ 2) | **24.0** (↑ 8) | **78.0** (↑ 9) | **52.0**(↑ 7) |

| | Place Object Stand | Place Phone Stand | Place Shoe | Put Object Cabinet | Rotate QRcode | Shake Bottle Horizontally | Shake Bottle |
|---|---|---|---|---|---|---|---|
| RDT | 15.0 | 15.0 | 35.0 | 33.0 | 50.0 | 84.0 | 74.0 |
| $\pi_{0.5}$ | 68.0 | 76.0 | 53.0 | 54.0 | 62.0 | **100.0** | 98.0 |
| $\pi_{0.5}$ + TACO (Ours) | **78.0** (↑ 10) | **86.0** (↑ 10) | **65.0** (↑ 12) | **56.0** (↑ 2) | **67.0** (↑ 5) | **100.0** (0) | **99.0** (↑ 1) |

| | Stack Blocks Three | Stack Blocks Two | Stack Bowls Three | Stack Bowls Two | Stamp Seal | Turn Switch | Average |
|---|---|---|---|---|---|---|---|
| RDT | 2.0 | 21.0 | 51.0 | 76.0 | 1.0 | 35.0 | 34.6 |
| $\pi_{0.5}$ | 43.0 | 81.0 | 64.0 | **93.0** | 26.0 | 42.0 | 59.3 |
| $\pi_{0.5}$ + TACO (Ours) | **45.0** (↑ 2) | **91.0** (↑ 10) | **68.0** (↑ 4) | **93.0** (0) | **38.0** (↑ 12) | **49.0** (↑ 7) | **64.0** (↑ 4.7) |

**Table 4** Evaluation on the Libero-long benchmark. Our method, TACO, is applied to both Pi0.5 and OpenVLA. For the autoregressive VLA architecture, we set $temperature = 1$ for action sampling. Results marked with ∗ are directly reported from Robomonkey (Kwok et al., 2025). *OpenVLA (reproduced)* denotes our own reproduction results, which serves as the baseline for our TACO implementation. We observe that TACO can be effectively applied to autoregressive VLA and consistently improves performance.

| | $\pi_{0.5}$ (Black et al., 2025) | + TACO (Ours) | OpenVLA* (Kim et al., 2024) | + Robomonkey* (Kwok et al., 2025) | OpenVLA (reproduced) | + TACO (Ours) |
|---|---|---|---|---|---|---|
| Soup and Sauce in Basket | 98.0 | **100.0**(↑ 2) | 36.0 | 59.0 | 60.0 | **66.0**(↑ 6) |
| Cheese and Butter in Basket | 100.0 | 96.0(↓ 4) | 70.0 | 79.0 | 76.0 | **82.0**(↑ 6) |
| Turn on Stove and Place Moka | 98.0 | 98.0(0) | **58.0** | **58.0** | **58.0** | 52.0(↓ 6) |
| Black Bowl in Drawer | 98.0 | **100.0**(↑ 2) | 36.0 | 37.0 | 36.0 | **50.0**(↑ 14) |
| Mugs on Plates | 98.0 | 98.0(0) | 42.0 | 55.0 | 32.0 | **50.0**(↑ 18) |
| Book in Caddy | 100.0 | 100.0(0) | 84.0 | 86.0 | 82.0 | **90.0**(↑ 8) |
| Mug and Pudding on Plate | **96.0** | 92.0(↓ 4) | 48.0 | 59.0 | **60.0** | 54.0(↓ 6) |
| Soup and Cheese in Basket | 94.0 | **100.0**(↑ 6) | 56.0 | 62.0 | 70.0 | **80.0**(↑ 10) |
| Moka Pots on Stove | 68.0 | **86.0**(↑ 16) | 26.0 | 26.0 | 20.0 | **28.0**(↑ 8) |
| Mug in Microwave | **98.0** | 96.0(↓ 2) | 42.0 | 44.0 | 46.0 | **48.0**(↑ 2) |
| Average | 94.8 | **96.6**(↑ 1.8) | 49.8 | 56.5 | 54.0 | **60.0**(↑ 6) |

unplugs the charger from the socket, and stores it back in the cabinet; and *Pick Books*, where the left and right arms each pick up and lift a book simultaneously.

These tasks cover a wide range of skills and scenarios, including human-robot interaction, dual-arm coordination, and long-horizon task execution, providing a comprehensive evaluation of the policy's general capabilities. For data collection, we perform teleoperation to record 100 episodes for each task. The collected data are then used to fine-tune the base policy $\pi_0$, upon which we deploy our proposed test-time scaling framework. In the real-world experiments, we simultaneously sample 30 action chunks at each decision step, with each action chunk having a length of 20.

### 5.2.2  Results

We address **Q2:** *What is the additional time cost introduced by our test-time scaling method? Can it run efficiently on real robots while maintaining performance gains?*

To intuitively demonstrate the time cost of parallel inference over multiple actions, we first compare the inference latency of two settings: (1) direct parallel inference of multiple actions, and (2) our optimized implementation that shares the observation $k, v$ cache across actions. As shown in Figure 4 (left), we evaluate the policy $\pi_0$ on a single RTX 4090 GPU. The results show that the $k, v$ caching optimization significantly accelerates inference, especially when the number of parallel actions is large. When sampling 32 actions simultaneously, the inference time is reduced by 73.2% compared to naïve multi-batch parallel inference.
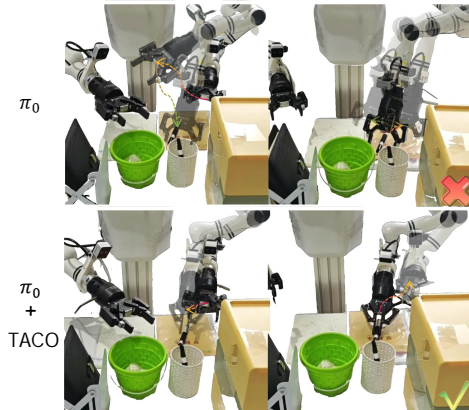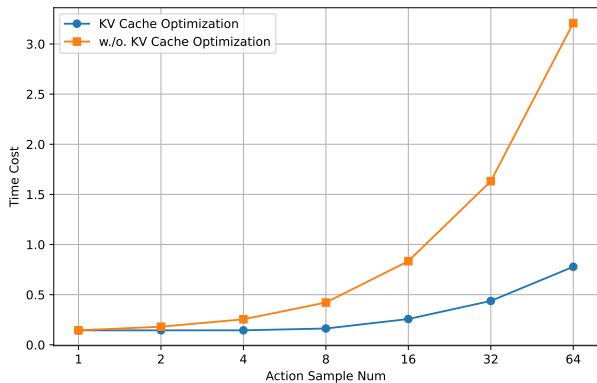
**Figure 4** (Left) Efficiency improvement from KV cache optimization. For each number of action samples, we repeated the inference 50 times and reported the average inference time. (Right) A key moment from the *Paper and Pen* task, where the robot needs to pick up a marker. The base policy $\pi_0$ often falls into a suboptimal mode, causing the arm to swing back and forth and hesitate before grasping, which leads to task failure. In contrast, our method enables $\pi_0$ to completely avoid such oscillations and hesitation during grasping.

**Table 5** We conducted an ablation study on the CFN network, the method for obtaining features, and the feature processing approach.

|  | Block Handover | Container Place | Diverse Bottles Pick | Average |
|---|---|---|---|---|
| $\pi_0$ **+ TACO (Ours)** | **62.0** | **40.0** | 27.0 | **43.0** |
| w./o. CFN | 52.0 | 32.0 | 26.0 | 36.7 |
| w./o. Feature Scaling | 51.0 | 37.0 | **30.0** | 39.3 |
| w./o. Internal Feature | 48.0 | 33.0 | 23.0 | 34.7 |

Therefore, our method remains efficient and low-latency when deployed on real robots.

The success rates in real-world experiments are shown in Figure 5. Our method achieves an average improvement of 16% over the base policy, with particularly notable gains in long-horizon tasks. Specifically, we observe a 25% improvement in the *Paper and Pen* task and a 15% improvement in the *Laptop* task. Based on our observations, the base policy typically fails in two ways: (1) imprecise grasping positions, and (2) suboptimal teleoperation data quality.

Our scaling framework improves action precision, consistent with the analysis observed in simulation experiments. More importantly, our approach effectively mitigates policy failures caused by imperfect teleoperation data. As shown in Figure 4 (right), during the *Pen Grasping* stage, the operator's varying grasp poses and timings lead to large state differences between grasping modes, resulting in a sparse observation–action distribution in the expert dataset. Consequently, the base policy often exhibits unstable behaviors such as failing to close the gripper properly, re-opening after grasping, or oscillating between two grasping modes. With our test-time scaling method, these issues are largely eliminated: the gripper consistently closes correctly, and the robot avoids suboptimal modes. We attribute this improvement to the filtering effect of our scaling framework, which suppresses suboptimal trajectories introduced by noise in the teleoperated dataset.
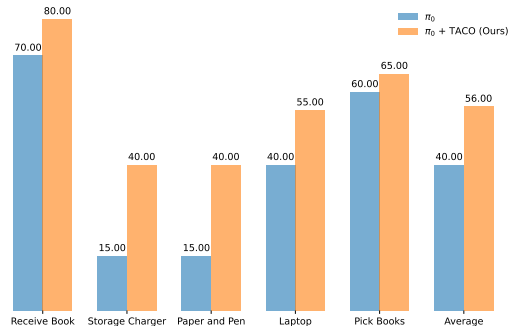


**Figure 5** Success rate (%) in real-world experiments. On our self-built *RealMan 75* dual-arm robotic platform, applying the proposed test-time scaling method improves both task accuracy and safety, resulting in an average 16% increase in success rate.

10

### 5.3 Ablation Studies and Analyses

To answer **Q3:** *How do design choices—such as the use of internal features and the CFN pseudo-count network—affect overall performance?* We designed a set of ablation experiments based on $\pi_0$ on the RobotWin benchmark. As shown in Table 5, the experiments include: (1) replacing CFN pseudo-counting of good features with a direct fit of the Euclidean distance between the predicted action from the features and the ground-truth action; (2) applying CFN for pseudo-counting without performing feature scaling on Internal Features; (3) using a CNN-based image encoder and an MLP-based action encoder to obtain features instead of relying on Internal Features.

Directly fitting the mapping from features to action errors increases the learning difficulty, as the model must capture not only optimal features but also suboptimal and poor ones. During test-time scaling, only the selection of the optimal action matters. Our experimental analysis suggests that the observed performance drop is caused by overfitting on the training set, leading to reduced generalization. When separate encoders are used instead of Internal Features and trained jointly with CFN, the resulting features tend to be highly similar, making accurate pseudo-counting difficult and causing a significant performance drop. This approach also reduces efficiency, as it fails to leverage VLA's inherent understanding of scenes and actions.

## 6    Conclusion

In this work, we introduce **TACO**, a test-time-scaling (TTS) framework that employs a lightweight pseudo-count estimator as a high-fidelity verifier for action chunks. TACO is compatible with various VLA models and consistently delivers substantial improvements over the base policy across a wide range of real-world and simulated experiments. Despite certain limitations—such as its inability to evaluate newly synthesized actions produced purely by perturbation and its reliance on the representational capability of the underlying VLA—our approach demonstrates that lightweight test-time scaling can effectively achieve anti-exploration and mitigate the instability and performance degradation caused by distribution shift in action prediction.

# References

Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. An optimistic perspective on offline reinforcement learning. In *International Conference on Machine Learning*, 2019.

Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-based offline reinforcement learning with diversified q-ensemble. In *Advances in Neural Information Processing Systems*, 2021. URL https://openreview.net/forum?id=ZUvaSolQZh3.

Johannes Anschütz and Arthur-César Le Bras. Prismatic dieudonné theory. In *Forum of Mathematics, Pi*, volume 11, page e2. Cambridge University Press, 2023.

Chenjia Bai, Lingxiao Wang, Zhuoran Yang, Zhi-Hong Deng, Animesh Garg, Peng Liu, and Zhaoran Wang. Pessimistic bootstrapping for uncertainty-driven offline reinforcement learning. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=Y4cs1Z3HnqL.

Chenjia Bai, Yang Zhang, Shuang Qiu, Qiaosheng Zhang, Kang Xu, and Xuelong Li. Online preference alignment for language models via count-based exploration. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=cfKZ5VrhXt.

Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos. Unifying count-based exploration and intrinsic motivation. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL https://proceedings.neurips.cc/paper_files/paper/2016/file/afda332245e2af431fb7b672a68b659d-Paper.pdf.

Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024.

Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding, Linxi Fan, Yu Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, et al. Gr00t n1: An open foundation model for generalist humanoid robots. *arXiv preprint arXiv:2503.14734*, 2025.

Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al. $\pi_0$: A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024.

Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail, Michael Robert Equi, Chelsea Finn, Niccolo Fusai, Manuel Y. Galliker, Dibya Ghosh, Lachy Groom, Karol Hausman, brian ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Devin LeBlanc, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Allen Z. Ren, Lucy Xiaoyang Shi, Laura Smith, Jost Tobias Springenberg, Kyle Stachowicz, James Tanner, Quan Vuong, Homer Walke, Anna Walling, Haohuan Wang, Lili Yu, and Ury Zhilinsky. $\pi_{0.5}$: a vision-language-action model with open-world generalization. In *9th Annual Conference on Robot Learning*, 2025. URL https://openreview.net/forum?id=vlhoswksBO.

Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.

Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*, 2024.

Qingwen Bu, Jisong Cai, Li Chen, Xiuqi Cui, Yan Ding, Siyuan Feng, Shenyuan Gao, Xindong He, Xuan Hu, Xu Huang, et al. Agibot world colosseo: A large-scale manipulation platform for scalable and intelligent embodied systems. *arXiv preprint arXiv:2503.06669*, 2025.

Remi Cadene, Simon Alibert, Alexander Soare, Quentin Gallouedec, Adil Zouitine, Steven Palma, Pepijn Kooijmans, Michel Aractingi, Mustafa Shukor, Dana Aubakirova, Martino Russi, Francesco Capuano, Caroline Pascal, Jade Choghari, Jess Moss, and Thomas Wolf. Lerobot: State-of-the-art machine learning for real-world robotics in pytorch. https://github.com/huggingface/lerobot, 2024.

Guoxin Chen, Minpeng Liao, Chengxi Li, and Kai Fan. Alphamath almost zero: Process supervision without process. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, *Advances in Neural Information Processing Systems*, volume 37, pages 27689–27724. Curran Associates, Inc., 2024. doi: 10.52202/079017-0870. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/30dfe47a3ccbee68cffa0c19ccb1bc00-Paper-Conference.pdf.

Tianxing Chen, Zanxin Chen, Baijun Chen, Zijian Cai, Yibin Liu, Zixuan Li, Qiwei Liang, Xianliang Lin, Yiheng Ge, Zhenyu Gu, et al. Robotwin 2.0: A scalable data generator and benchmark with strong domain randomization for robust bimanual robotic manipulation. *arXiv preprint arXiv:2506.18088*, 2025.

Xi Chen, Josip Djolonga, Piotr Padlewski, Basil Mustafa, Soravit Changpinyo, Jialin Wu, Carlos Riquelme Ruiz, Sebastian Goodman, Xiao Wang, Yi Tay, et al. Pali-x: On scaling up a multilingual vision and language model. *arXiv preprint arXiv:2305.18565*, 2023.

Open X-Embodiment Collaboration. Open X-Embodiment: Robotic learning datasets and RT-X models. https://arxiv.org/abs/2310.08864, 2023.

Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. In *International Conference on Machine Learning*, pages 8469–8488. PMLR, 2023.

Max Du and Shuran Song. Dynaguide: Steering diffusion polices with active dynamic guidance. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. URL https://openreview.net/forum?id=XOw7Yf8qN3.

Ryan Ehrlich, Bradley Brown, Jordan Juravsky, Ronald Clark, Christopher Ré, and Azalia Mirhoseini. Codemonkeys: Scaling test-time compute for software engineering. *arXiv preprint arXiv:2501.14723*, 2025.

Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning. *ArXiv*, abs/2106.06860, 2021. URL https://api.semanticscholar.org/CorpusID:235422620.

Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In *International conference on machine learning*, pages 2052–2062. PMLR, 2019.

Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019.

Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with discrete world models. *arXiv preprint arXiv:2010.02193*, 2020.

Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

Zhi Hou, Tianyi Zhang, Yuwen Xiong, Haonan Duan, Hengjun Pu, Ronglei Tong, Chengyang Zhao, Xizhou Zhu, Yu Qiao, Jifeng Dai, et al. Dita: Scaling diffusion transformer for generalist vision-language-action policy. *arXiv preprint arXiv:2503.19757*, 2025.

Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Frederik Ebert, Corey Lynch, Sergey Levine, and Chelsea Finn. BC-z: Zero-shot task generalization with robotic imitation learning. In *5th Annual Conference on Robot Learning*, 2021. URL https://openreview.net/forum?id=8kbp23tSGYv.

Suhyeok Jang, Dongyoung Kim, Changyeon Kim, Youngsuk Kim, and Jinwoo Shin. Verifier-free test-time sampling for vision language action models. *ArXiv*, abs/2510.05681, 2025. URL https://api.semanticscholar.org/CorpusID:281886270.

Zhi Jing, Siyuan Yang, Jicong Ao, Ting Xiao, Yugang Jiang, and Chenjia Bai. Humanoidgen: Data generation for bimanual dexterous manipulation via llm reasoning. *arXiv preprint arXiv:2507.00833*, 2025.

Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis, et al. Droid: A large-scale in-the-wild robot manipulation dataset. *arXiv preprint arXiv:2403.12945*, 2024.

Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-based offline reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 33, pages 21810–21823. Curran Associates, Inc., 2020.

Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan P Foster, Pannag R Sanketi, Quan Vuong, Thomas Kollar, Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. OpenVLA: An open-source vision-language-action model. In *8th Annual Conference on Robot Learning*, 2024. URL https://openreview.net/forum?id=ZMnD6QZAE6.

Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-learning. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=68n2s9ZJWF8.

Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. *Advances in neural information processing systems*, 32, 2019.

Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. *Advances in neural information processing systems*, 33:1179–1191, 2020.

Jacky Kwok, Christopher Agia, Rohan Sinha, Matt Foutter, Shulu Li, Ion Stoica, Azalia Mirhoseini, and Marco Pavone. Robomonkey: Scaling test-time sampling and verification for vision-language-action models. In Joseph Lim, Shuran Song, and Hae-Won Park, editors, *Proceedings of The 9th Conference on Robot Learning*, volume 305 of *Proceedings of Machine Learning Research*, pages 3200–3217. PMLR, 27–30 Sep 2025. URL https://proceedings.mlr.press/v305/kwok25a.html.

Sergey Levine, Aviral Kumar, G. Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *ArXiv*, abs/2005.01643, 2020a. URL https://api.semanticscholar.org/CorpusID:218486979.

Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020b.

Xinghang Li, Peiyan Li, Minghuan Liu, Dong Wang, Jirong Liu, Bingyi Kang, Xiao Ma, Tao Kong, Hanbo Zhang, and Huaping Liu. Towards generalist robot policies: What matters in building vision-language-action models. *arXiv preprint arXiv:2412.14058*, 2024a.

Xuanlin Li, Kyle Hsu, Jiayuan Gu, Karl Pertsch, Oier Mees, Homer Rich Walke, Chuyuan Fu, Ishikaa Lunawat, Isabel Sieh, Sean Kirmani, Sergey Levine, Jiajun Wu, Chelsea Finn, Hao Su, Quan Vuong, and Ted Xiao. Evaluating real-world robot manipulation policies in simulation. *arXiv preprint arXiv:2405.05941*, 2024b.

Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=PqvMRDCJT9t.

Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky TQ Chen, David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow matching guide and code. *arXiv preprint arXiv:2412.06264*, 2024.

Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. *arXiv preprint arXiv:2306.03310*, 2023a.

Songming Liu, Lingxuan Wu, Bangguo Li, Hengkai Tan, Huayu Chen, Zhengyi Wang, Ke Xu, Hang Su, and Jun Zhu. RDT-1b: a diffusion foundation model for bimanual manipulation. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=yAzN4tz7oI.

Xingchao Liu, Chengyue Gong, and qiang liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *The Eleventh International Conference on Learning Representations*, 2023b. URL https://openreview.net/forum?id=XVjTT1nw5z.

Sam Lobel, Akhil Bagaria, and George Konidaris. Flipping coins to estimate pseudocounts for exploration in reinforcement learning. In *International Conference on Machine Learning*, pages 22594–22613. PMLR, 2023.

Yao Mu, Tianxing Chen, Zanxin Chen, Shijia Peng, Zhiqian Lan, Zeyu Gao, Zhixuan Liang, Qiaojun Yu, Yude Zou, Mingkun Xu, et al. Robotwin: Dual-arm robot benchmark with generative digital twins. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 27649–27660, 2025.

Mitsuhiko Nakamoto, Oier Mees, Aviral Kumar, and Sergey Levine. Steering your generalists: Improving robotic foundation models via value guidance. In *8th Annual Conference on Robot Learning*, 2024. URL https://openreview.net/forum?id=6GlpzC9Po.

Alexander Nikulin, Vladislav Kurenkov, Denis Tarasov, and Sergey Kolesnikov. Anti-exploration by random network distillation. In *International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 26116–26130, 2023.

Georg Ostrovski, Marc G Bellemare, Aäron van den Oord, and Rémi Munos. Count-based exploration with neural density models. In *International Conference on Machine Learning*. PMLR, 2017.

Abby O'Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, Albert Tung, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anchit Gupta, Andrew Wang, Anikait Singh, Animesh Garg, Aniruddha Kembhavi, Annie Xie, Anthony Brohan, Antonin Raffin, Archit Sharma, Arefeh Yavary, Arhan Jain, Ashwin Balakrishna, Ayzaan

Wahid, Ben Burgess-Limerick, Beomjoon Kim, Bernhard Schölkopf, Blake Wulfe, Brian Ichter, Cewu Lu, Charles Xu, Charlotte Le, Chelsea Finn, Chen Wang, Chenfeng Xu, Cheng Chi, Chenguang Huang, Christine Chan, Christopher Agia, Chuer Pan, Chuyuan Fu, Coline Devin, Danfei Xu, Daniel Morton, Danny Driess, Daphne Chen, Deepak Pathak, Dhruv Shah, Dieter Büchler, Dinesh Jayaraman, Dmitry Kalashnikov, Dorsa Sadigh, Edward Johns, Ethan Foster, Fangchen Liu, Federico Ceola, Fei Xia, Feiyu Zhao, Freek Stulp, Gaoyue Zhou, Gaurav S. Sukhatme, Gautam Salhotra, Ge Yan, Gilbert Feng, Giulio Schiavi, Glen Berseth, Gregory Kahn, Guanzhi Wang, Hao Su, Hao-Shu Fang, Haochen Shi, Henghui Bao, Heni Ben Amor, Henrik I Christensen, Hiroki Furuta, Homer Walke, Hongjie Fang, Huy Ha, Igor Mordatch, Ilija Radosavovic, Isabel Leal, Jacky Liang, Jad Abou-Chakra, Jaehyung Kim, Jaimyn Drake, Jan Peters, Jan Schneider, Jasmine Hsu, Jeannette Bohg, Jeffrey Bingham, Jeffrey Wu, Jensen Gao, Jiaheng Hu, Jiajun Wu, Jialin Wu, Jiankai Sun, Jianlan Luo, Jiayuan Gu, Jie Tan, Jihoon Oh, Jimmy Wu, Jingpei Lu, Jingyun Yang, Jitendra Malik, João Silvério, Joey Hejna, Jonathan Booher, Jonathan Tompson, Jonathan Yang, Jordi Salvador, Joseph J. Lim, Junhyek Han, Kaiyuan Wang, Kanishka Rao, Karl Pertsch, Karol Hausman, Keegan Go, Keerthana Gopalakrishnan, Ken Goldberg, Kendra Byrne, Kenneth Oslund, Kento Kawaharazuka, Kevin Black, Kevin Lin, Kevin Zhang, Kiana Ehsani, Kiran Lekkala, Kirsty Ellis, Krishan Rana, Krishnan Srinivasan, Kuan Fang, Kunal Pratap Singh, Kuo-Hao Zeng, Kyle Hatch, Kyle Hsu, Laurent Itti, Lawrence Yunliang Chen, Lerrel Pinto, Li Fei-Fei, Liam Tan, Linxi Jim Fan, Lionel Ott, Lisa Lee, Luca Weihs, Magnum Chen, Marion Lepert, Marius Memmel, Masayoshi Tomizuka, Masha Itkina, Mateo Guaman Castro, Max Spero, Maximilian Du, Michael Ahn, Michael C. Yip, Mingtong Zhang, Mingyu Ding, Minho Heo, Mohan Kumar Srirama, Mohit Sharma, Moo Jin Kim, Naoaki Kanazawa, Nicklas Hansen, Nicolas Heess, Nikhil J Joshi, Niko Suenderhauf, Ning Liu, Norman Di Palo, Nur Muhammad Mahi Shafiullah, Oier Mees, Oliver Kroemer, Osbert Bastani, Pannag R Sanketi, Patrick Tree Miller, Patrick Yin, Paul Wohlhart, Peng Xu, Peter David Fagan, Peter Mitrano, Pierre Sermanet, Pieter Abbeel, Priya Sundaresan, Qiuyu Chen, Quan Vuong, Rafael Rafailov, Ran Tian, Ria Doshi, Roberto Martín-Martín, Rohan Baijal, Rosario Scalise, Rose Hendrix, Roy Lin, Runjia Qian, Ruohan Zhang, Russell Mendonca, Rutav Shah, Ryan Hoque, Ryan Julian, Samuel Bustamante, Sean Kirmani, Sergey Levine, Shan Lin, Sherry Moore, Shikhar Bahl, Shivin Dass, Shubham Sonawani, Shuran Song, Sichun Xu, Siddhant Haldar, Siddharth Karamcheti, Simeon Adebola, Simon Guist, Soroush Nasiriany, Stefan Schaal, Stefan Welker, Stephen Tian, Subramanian Ramamoorthy, Sudeep Dasari, Suneel Belkhale, Sungjae Park, Suraj Nair, Suvir Mirchandani, Takayuki Osa, Tanmay Gupta, Tatsuya Harada, Tatsuya Matsushima, Ted Xiao, Thomas Kollar, Tianhe Yu, Tianli Ding, Todor Davchev, Tony Z. Zhao, Travis Armstrong, Trevor Darrell, Trinity Chung, Vidhi Jain, Vincent Vanhoucke, Wei Zhan, Wenxuan Zhou, Wolfram Burgard, Xi Chen, Xiaolong Wang, Xinghao Zhu, Xinyang Geng, Xiyuan Liu, Xu Liangwei, Xuanlin Li, Yao Lu, Yecheng Jason Ma, Yejin Kim, Yevgen Chebotar, Yifan Zhou, Yifeng Zhu, Yilin Wu, Ying Xu, Yixuan Wang, Yonatan Bisk, Yoonyoung Cho, Youngwoon Lee, Yuchen Cui, Yue Cao, Yueh-Hua Wu, Yujin Tang, Yuke Zhu, Yunchu Zhang, Yunfan Jiang, Yunshuang Li, Yunzhu Li, Yusuke Iwasawa, Yutaka Matsuo, Zehan Ma, Zhuo Xu, Zichen Jeff Cui, Zichen Zhang, and Zipeng Lin. Open x-embodiment: Robotic learning datasets and rt-x models : Open x-embodiment collaboration0. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6892–6903, 2024. doi: 10.1109/ICRA57147.2024.10611477.

Delin Qu, Haoming Song, Qizhi Chen, Yuanqi Yao, Xinyi Ye, Yan Ding, Zhigang Wang, JiaYuan Gu, Bin Zhao, Dong Wang, et al. Spatialvla: Exploring spatial representations for visual-language-action model. *arXiv preprint arXiv:2501.15830*, 2025.

Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart J. Russell. Bridging offline reinforcement learning and imitation learning: A tale of pessimism. *IEEE Transactions on Information Theory*, 68:8156–8196, 2021. URL https://api.semanticscholar.org/CorpusID:232307477.

Shideh Rezaeifar, Robert Dadashi, Nino Vieillard, Léonard Hussenot, Olivier Bachem, Olivier Pietquin, and Matthieu Geist. Offline reinforcement learning as anti-exploration. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 8106–8114, 2022.

Mustafa Shukor, Dana Aubakirova, Francesco Capuano, Pepijn Kooijmans, Steven Palma, Adil Zouitine, Michel Aractingi, Caroline Pascal, Martino Russi, Andres Marafioti, et al. Smolvla: A vision-language-action model for affordable and efficient robotics. *arXiv preprint arXiv:2506.01844*, 2025.

Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. pmlr, 2015.

Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=PxTIG12RRHS.

Richard S. Sutton, David A. McAllester, Satinder Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Neural Information Processing Systems*, 1999. URL https://api.semanticscholar.org/CorpusID:1211821.

Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023. URL https://arxiv.org/abs/2307.09288.

Homer Rich Walke, Kevin Black, Tony Z. Zhao, Quan Vuong, Chongyi Zheng, Philippe Hansen-Estruch, Andre Wang He, Vivek Myers, Moo Jin Kim, Max Du, Abraham Lee, Kuan Fang, Chelsea Finn, and Sergey Levine. Bridgedata v2: A dataset for robot learning at scale. In *7th Annual Conference on Robot Learning*, 2023. URL https://openreview.net/forum?id=f55MlAT1Lu.

Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.

Junjie Wen, Yichen Zhu, Jinming Li, Zhibin Tang, Chaomin Shen, and Feifei Feng. Dexvla: Vision-language model with plug-in diffusion expert for general robot control. *arXiv preprint arXiv:2502.05855*, 2025.

Kun Wu, Chengkai Hou, Jiaming Liu, Zhengping Che, Xiaozhu Ju, Zhuqin Yang, Meng Li, Yinuo Zhao, Zhiyuan Xu, Guang Yang, et al. Robomind: Benchmark on multi-embodiment intelligence normative data for robot manipulation. *arXiv preprint arXiv:2412.13877*, 2024.

Yifan Wu, G. Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. *ArXiv*, abs/1911.11361, 2019. URL https://api.semanticscholar.org/CorpusID:208291277.

Tianpei Yang, Hongyao Tang, Chenjia Bai, Jinyi Liu, Jianye Hao, Zhaopeng Meng, Peng Liu, and Zhen Wang. Exploration in deep reinforcement learning: From single-agent to multiagent domain. *IEEE Transactions on Neural Networks and Learning Systems*, 35:8762–8782, 2021. URL https://api.semanticscholar.org/CorpusID:250493352.

Yang Zhang, Chenjia Bai, Bin Zhao, Junchi Yan, Xiu Li, and Xuelong Li. Decentralized transformers with centralized aggregation are sample-efficient multi-agent world models. *Transactions on Machine Learning Research*, 2025a. ISSN 2835-8856. URL https://openreview.net/forum?id=xT8BEgXmVc.

Yang Zhang, Xinran Li, Jianing Ye, Shuang Qiu, Delin Qu, Xiu Li, Chongjie Zhang, and Chenjia Bai. Revisiting multi-agent world modeling from a diffusion-inspired perspective. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025b. URL https://openreview.net/forum?id=rRxFIOoEeF.

Yang Zhang, Chenwei Wang, Ouyang Lu, Yuan Zhao, Yunfei Ge, Zhenglong Sun, Xiu Li, Chi Zhang, Chenjia Bai, and Xuelong Li. Align-then-steer: Adapting the vision-language action models through unified latent guidance. *arXiv preprint arXiv:2509.02055*, 2025c.

Yang Zhang, Shixin Yang, Chenjia Bai, Fei Wu, Xiu Li, Zhen Wang, and Xuelong Li. Towards efficient llm grounding for embodied multi-agent collaboration. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 1663–1699, 2025d.

Qingqing Zhao, Yao Lu, Moo Jin Kim, Zipeng Fu, Zhuoyang Zhang, Yecheng Wu, Zhaoshuo Li, Qianli Ma, Song Han, Chelsea Finn, Ankur Handa, Ming-Yu Liu, Donglai Xiang, Gordon Wetzstein, and Tsung-Yi Lin. Cot-vla: Visual chain-of-thought reasoning for vision-language-action models. *2025 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1702–1713, 2025. URL https://api.semanticscholar.org/CorpusID:277435005.

Yifan Zhong, Xuchuan Huang, Ruochong Li, Ceyao Zhang, Zhang Chen, Tianrui Guan, Fanlian Zeng, Ka Num Lui, Yuyao Ye, Yitao Liang, Yaodong Yang, and Yuanpei Chen. Dexgraspvla: A vision-language-action framework towards general dexterous grasping, 2025. URL https://arxiv.org/abs/2502.20900.

Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, Ayzaan Wahid, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *Conference on Robot Learning*, pages 2165–2183. PMLR, 2023.

# A   Algorithmic Description of TACO during Inference

As shown in Algorithm 1, we provide an algorithmic description for the action-filtering procedure during inference time, illustrating how the CFN verifier leverages internal representations to obtain pseudo-counts, as well as the workflow of the KV Cache Optimization. Here, we take flow-based VLAs as representative examples.

---

**Algorithm 1:** Anti-Exploration Sampling for TTS

---

**Input:**  Observation $\mathbf{o}_t$; instruction $l$;
VLA policy $\pi_\theta$; CFN verifier $f_\phi$;
Number of candidates $M$.
**Output:** Selected action chunk $\mathbf{a}^*_{t:t+H}$

Compute context KV cache $\mathcal{C}_{\mathrm{KV}} = \mathrm{ContextForward}(\mathbf{o}_t, l)$ in $\pi_\theta$;
Sample $M$ noise candidates $\{\epsilon_i\}_{i=1}^M \sim \mathcal{N}(0, I)$;
**for** $i = 1$ **to** $M$ **do**
> Compute the predicted clean action chunk and last hidden state
> $(\hat{\mathbf{a}}^{(i)}_{t:t+H}, h^{(i)}_\theta) = \pi_\theta(\epsilon_i, \mathcal{C}_{\mathrm{KV}})$, i.e., $\pi_\theta(\mathbf{o}_t, l; \epsilon_i)$;

**for** $i = 1$ **to** $M$ **do**
> Compute the pseudo-count of the generated action chunk $\hat{N}_{\mathcal{D}_{\mathrm{sft}}}(\mathbf{o}_t, l, \hat{\mathbf{a}}^{(i)}_{t:t+H}) = \frac{1}{\|f_\phi(h^{(i)}_\theta)\|^2}$;

Select $i^* = \arg\max_i \hat{N}_{\mathcal{D}_{\mathrm{sft}}}(\mathbf{o}_t, l, \hat{\mathbf{a}}^{(i)}_{t:t+H})$;
**return** $\mathbf{a}^*_{t:t+H} = \hat{\mathbf{a}}^{(i^*)}_{t:t+H}$

---

# B   Algorithmic Description of TACO during Training

Algorithm 2 summarizes how we train the coupled pseudo-count estimator on top of a VLA. The goal is to efficiently obtain in-distribution, high-fidelity features by repeatedly querying the VLA under different noise conditions, then selecting the representation that best aligns with the clean action. Training a lightweight CFN on top of these coupled features yields a pseudo-count estimator that reliably reflects how close a candidate action lies to the support of the dataset $\mathcal{D}_{\mathrm{sft}}$, enabling effective anti-exploration during inference. Here, we also take flow-based VLAs as representative examples.

# C   Additional Preliminary

**Flow Matching.** Flow matching (Liu et al., 2023b; Lipman et al., 2023, 2024) is a generative modeling framework that learns a continuous velocity field to transform a simple prior distribution into a target data distribution. Similar to diffusion models, it defines a time-dependent transformation; however, the dynamics are governed by an ordinary differential equation (ODE) rather than a stochastic process.

Formally, let $p(x)$ denote the target distribution on $\mathbb{R}^k$, and let $v(u, x) : [0, 1] \times \mathbb{R}^k \to \mathbb{R}^k$ be a time-dependent velocity field. Its flow $\psi(u, x)$ is defined by

$$\frac{\mathrm{d}}{\mathrm{d}u}\psi(u, x) = v(u, \psi(u, x)), \quad \psi(0, x) = x. \tag{10}$$

By integrating this ODE from $u = 0$ to $u = 1$, samples from a simple prior (e.g., $\mathcal{N}(0, I)$) are continuously mapped to the target distribution.

Flow matching provides a stable and efficient alternative to diffusion-based generation and offers a principled continuous-time formulation that can be extended to policy learning or trajectory generation in robotic systems.

---

**Algorithm 2:** Coupled Pseudo-Count Estimator Training for VLAs

---

**Input:** SFT dataset $\mathcal{D}_{\text{sft}} = \{(\mathbf{o}, l, \mathbf{a})\}$;
VLA policy $\pi_\theta$; CFN head $f_\phi$; learning rate $\eta$;
Noise schedule $\{\sigma_i\}_{i=1}^N$; number of queries $N$;
Number of training steps $T$; batch size $B$.
**Output:** Trained CFN $f_\phi$; high-fidelity feature set $\mathcal{D}_h$

// High-Fidelity Feature Search
Initialize feature set: $\mathcal{D}_h \leftarrow \emptyset$;
**foreach** $(\mathbf{o}, l, \mathbf{a}) \in \mathcal{D}_{sft}$ **do**
    **for** $i = 1$ **to** $N$ **do**
        Corrupt the clean action $\mathbf{a}_{\sigma_i} = \sigma_i \mathbf{a} + (1 - \sigma_t)\epsilon$ with $\epsilon \sim \mathcal{N}(0, I)$;
        Compute the predicted clean action chunk and last hidden state $(\mathbf{a}_{\text{pre}}^{(i)}, h_\theta^{(i)}) = \pi_\theta(\mathbf{o}, l, \mathbf{a}_{\sigma_i})$;
        Calculate the error between the prediction and the ground truth label $e_i \leftarrow \|\mathbf{a}_{\text{pre}}^{(i)} - \mathbf{a}\|_2$;
    Select $i^\star = \arg\min_i e_i$;
    Add the internal feature of the closest prediction into the feature set $\mathcal{D}_h \leftarrow \mathcal{D}_h \cup \{h_\theta^{(i^\star)}\}$;

// CFN Training on Coupled Features;
**for** $t = 1$ **to** $T$ **do**
    Sample mini-batch $\{h_k\}_{k=1}^B \sim \mathcal{D}_h$;
    Compute $\mathcal{L}_{\text{CFN}}$ according to Eq. (1);
    Update CFN parameters: $\phi \leftarrow \phi - \eta \nabla_\phi \mathcal{L}_{\text{CFN}}$;
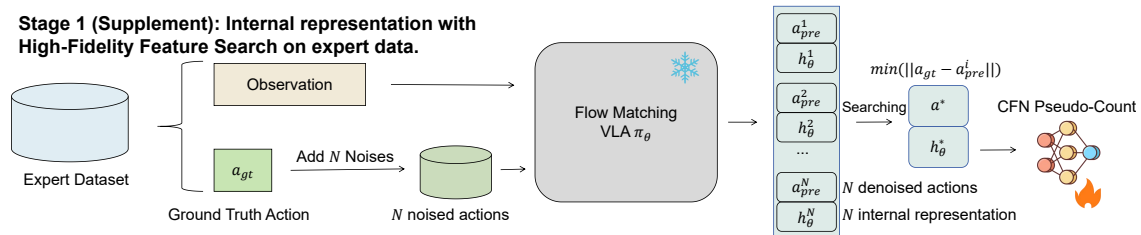**return** $f_\phi$, $\mathcal{D}_h$

---



**Figure 6** During training, each internal representation undergoes a search process that selects the best one from $N$ candidates to maximally preserve the information of the ground-truth action.

## D    Training Implementation Details

**Internal Representation Details.** For the two policies $\pi_0$ and $\pi_{0.5}$, we build the training and deployment pipeline based on the *Lerobot* framework (Cadene et al., 2024). We first modify the model's action-output function so that it simultaneously returns both the action and the internal representation. Concretely, we extract the first action token from the final hidden layer. Since the final hidden layer encodes an abstract representation of the entire input sequence—and, through the self-attention mechanism, the last token's vector contains contextual information of the whole sequence—we treat this vector as a compact representation of the input. Because the action expert used in both $\pi_0$ and $\pi_{0.5}$ adopts a bidirectional-attention Transformer architecture, in practice we take the first token of the final hidden layer inside the action expert as the internal representation. Its dimensionality is 1024. For *OpenVLA*, we instead take the last token of its final hidden layer as the internal representation, whose dimensionality is 4096. We set the temperature to 1 and leverage the inherent randomness in decoding to sample diverse actions and features.

**Base Policy and Representation Dataset.** Taking *Robotwin1.0* (Mu et al., 2025) as an example, we first collect 50 trajectories for each task using the provided scripted policies to construct the dataset. Expert data collected via scripted policies can typically avoid a large amount of human-induced noise, enabling the construction of high-quality datasets at low cost (Jing et al., 2025). We then adopt the $\pi_0$ base model released by LeRobot

on Hugging Face, and perform full-parameter fine-tuning on all task data using an NVIDIA H100 with a batch size of 48 for 30k steps, resulting in the base policy. We inject 10% noise into the ground-truth actions and apply a single denoising step using the $\pi_0$ action expert to obtain the internal representations. As shown in Figure 6, for each triplet $(\mathbf{o}, l, \mathbf{a})$, we additionally sample $N = 50$ noise instances to generate: $\{(\mathbf{a}_{\text{pre}}^{(i)}, h_\theta^{(i)})\}_{i=1}^N$, and select the sample $\mathbf{a}_{\text{pre}}^{(*)}$ that is closest to the ground-truth action, which is then stored as input for subsequent training.

For *Robotwin2.0* (Chen et al., 2025), we similarly use the LeRobot $\pi_0$ base model from Hugging Face and fine-tune it on an NVIDIA H100 with a batch size of 32 for 10k steps across all tasks. When obtaining internal representations, we empirically set the noise level to 100% relative to the ground-truth actions. All other settings remain unchanged.

For *Libero* (Liu et al., 2023a), we directly use the official LeRobot model `pi0_libero_finetuned` as well as the open-source `openvla-7b-finetuned-libero-10` available on Hugging Face as our base policies.

For *Simpler* (Li et al., 2024b), we use the *Bridge V2* dataset and perform full-parameter fine-tuning of $\pi_0$ on 8 NVIDIA H100 GPUs with a batch size of 256, which we take as the base policy.

**Training Parameters** We take the training parameters on RobotWin2.0 as an example, and the remaining training configurations are largely similar. The training process of the CFN follows the settings listed in Table 6. The Adam optimizer is selected for parameter updates due to its efficiency in handling sparse gradients; The OneCycleLR scheduler is employed to dynamically adjust the learning rate: it linearly increases the learning rate from the initial value ($10^{-4}$) to the maximum value ($10^{-3}$) in the first half of training, and then cosine-anneals the learning rate to a lower value in the second half, which helps to balance convergence speed and generalization performance; Gradient accumulation is adopted with a step size of 2, i.e., gradients from 2 consecutive batches are accumulated before updating the model parameters. This effectively simulates a larger batch size while avoiding excessive GPU memory consumption.

**Table 6** Training parameters of CFN

| Parameter Name | Value |
| --- | --- |
| Batch Size | 512 |
| Action Chunk Size | 50 |
| CFN Output Dimension (`cfn_output_dim`) | 20 |
| CFN Hidden Dimension (`cfn_hidden_dim`) | 1536 |
| Total Training Epochs (`total_epochs`) | 16 |
| Optimizer | Adam |
| Initial Learning Rate | $10^{-4}$ |
| Learning Rate Scheduler | OneCycleLR |
| Maximum Learning Rate of Scheduler | $10^{-3}$ |
| Annealing Strategy | Cosine Annealing |
| Gradient Accumulation Steps (`grad_accum_steps`) | 2 |

# E   Training Procedure Details

## E.1   Network Architecture of CFN

The Coin Flipping Network (CFN) adopted in this work is a lightweight multi-layer perceptron (MLP), designed for efficient feature processing and target prediction. Its overall structure consists of four core components: a feature scaling module, an input projection layer, a stack of one or more MLP blocks, and an output projection layer. Detailed implementations are as follows:

**Feature Scaling** Before entering the network, all input features are scaled by a constant factor (set to 10 in this work) to enlarge the inter-feature distances. This operation helps reduce the interference between features during subsequent processing and improves the stability of feature discrimination.

**MLP Block** Each MLP block serves as the basic computational unit, employing a residual connection mechanism to mitigate the vanishing gradient problem during training. The structure of a single MLP block includes:

- Two fully connected (FC) layers, where the first layer maps the input feature dimension to 4 times the hidden dimension, and the second layer projects it back to the original hidden dimension;

- The GELU (Gaussian Error Linear Units) activation function is used to introduce non-linearity;

- A Dropout layer for regularization to avoid overfitting;

- A LayerNorm layer applied after the residual connection to stabilize the training process.

Weight initialization follows the Xavier uniform distribution, and biases are initialized to zero.

The specific workflow of the CFN network is as follows: *Feature Scaling*: All input features are multiplied by 10 to broaden their distribution range and improve feature separability; *Input Projection*: The scaled features are projected into the predefined hidden dimension using a fully connected (FC) layer, followed by GELU activation and LayerNorm normalization; *Feature Enhancement*: The normalized features are then passed through an MLP block for deep feature extraction and enhancement; *Output Projection*: Finally, the enhanced features are mapped to the target output dimension via the last FC layer.

## E.2  CFN Loss.

In our setting, the CFN is trained not on raw states but on the high-fidelity feature set $\mathcal{D}_h = \{h_\theta^{(i^\star)}\}$ extracted from the frozen VLA backbone. Following the same formulation as Eq. (1), the learning objective becomes

$$\mathcal{L}_{\text{CFN}}(\phi) = \mathbb{E}_{(h_\theta, \mathbf{c}) \sim \mathcal{D}_h} \left[ \| f_\phi(h_\theta) - \mathbf{c} \|^2 \right], \tag{11}$$

where each feature $h_\theta$ is paired with a randomly sampled Rademacher vector $\mathbf{c} \sim \{-1, 1\}^d$. Minimizing this distance-based loss encourages $f_\phi$ to approximate the expected coin-flip outcome associated with each feature embedding. As in the original CFN formulation, the norm of the optimal solution, $\| f_\phi(h_\theta) \|^2$, implicitly encodes the inverse visitation frequency of $h_\theta$. Hence, Eq. (9) can be directly interpreted as a distance-derived pseudo-count: features that are frequently observed during training yield smaller reconstruction errors and larger norms, whereas novel or out-of-support features correspond to larger distances in the CFN output space and smaller pseudo-count estimates. This formulation allows the CFN to act as a smooth verifier over the feature manifold without requiring explicit density estimation.

## E.3  Random Prior Initialization

To implement a low-count initialization for unseen features, we decompose the CFN output into a trainable component and a frozen random prior:

$$f_\phi(h) = \hat{f}_\phi(h) + f_{\text{prior}}(h), \tag{12}$$

where $f_{\text{prior}}$ is a randomly initialized, non-trainable network. We normalize the prior using running statistics so that each output dimension satisfies

$$\mathbb{E}_{h \sim \mathcal{D}_h} \left[ f_{\text{prior}}^{(i)}(h)^2 \right] = 1, \qquad \forall i \in \{1, \ldots, d\}, \tag{13}$$

which yields the expected squared norm

$$\mathbb{E}_{h \sim \mathcal{D}_h} \left[ \| f_{\text{prior}}(h) \|^2 \right] = d. \tag{14}$$

For features $h$ in regions where the learned component has not yet adapted, we have $\hat{f}_\phi(h) \approx 0$, and thus

$$\| f_\phi(h) \|^2 \approx \| f_{\text{prior}}(h) \|^2 \approx d. \tag{15}$$

Since our pseudo-count is defined up to a proportionality constant, we calibrate it such that

$$\hat{N}_{\mathcal{D}_{\text{sft}}}(h) \approx \frac{d}{\| f_\phi(h) \|^2} \approx 1. \tag{16}$$

In this way, unseen regions of the representation space are initialized with a pseudo-count of 1, and the influence of the random prior naturally diminishes as $\hat{f}_\phi$ is updated during training.

**Figure 7** Illustration of the benchmarks.

# F    Additional Simulation Experiment Details

**Benchmarks.** The schematic illustration of each benchmark task is shown in Figure 7. *RoboTwin 1.0* (Mu et al., 2025): One of the earlier dual-arm robot manipulation benchmarks based on generative digital twins, providing both real-world and simulated demonstration data for dual-arm tasks. *RoboTwin 2.0* (Chen et al., 2025): A scalable extension of RoboTwin that supports automatic large-scale data generation, strong domain randomization across clutter, lighting, background, tabletop height and language, and evaluation across multiple robot embodiments for robust generalization. *LIBERO* (Liu et al., 2023a): A lifelong robot learning benchmark designed to study knowledge transfer in related objects/skills, comprising 130 manipulation tasks across four task suites and supporting research in multi-task and lifelong learning for robotics. *SIMPLER* (Li et al., 2024b): A simulation-based evaluation environment set (Simulated Manipulation Policy Evaluation for Real Robot Setups) built on SAPIEN + ManiSkill2, aimed at providing scalable, reproducible simulation evaluations of manipulation policies that correlate well with real-world performance.

**Further Experimental Exploration on Representations.** We conducted extensive experiments using various methods for obtaining representations, including employing CNN and MLP as encoders (jointly optimized with the CFN parameters), utilizing internal representations, and combining internal representations with High-Fidelity Feature Search. As shown in Figure 8, we visualize how 1/Count evolves during the denoising process across 50 different initial noise configurations in the Block Handover task under different representation settings.

When using CNN and MLP as encoders, the pseudo-counter gradually fails to distinguish between different noise samples as the number of denoising steps increases. Eventually, the values converge to almost a single point, and the counter even reaches its maximum value before complete denoising, indicating that the counter is overly tolerant to actions.

In contrast, when using internal representations, the counter maintains a high degree of discrimination even after the noise has been completely removed. However, it shows limited distinction between better and worse noise samples. We attribute this to the fact that noise injection during the internal representation extraction process partially corrupts the ground-truth action information.

Finally, when internal representations are combined with High-Fidelity Feature Search, the search procedure effectively preserves the information of the ground-truth action by exploring the predicted action space. As a result, the model successfully distinguishes between better and worse noise samples once denoising is completed.

**More experiments on the choice of pseudo-counter.** To demonstrate the role and generalization ability of CFN, we compare it with RND (Nikulin et al., 2023) as a replacement pseudo-counter on a subset of tasks. RND employs a fixed random target network and a learned predictor whose estimation error serves as an intrinsic reward, enabling a lightweight and robust measure of feature novelty. We use this novelty as the pseudo-count value, where higher novelty corresponds to a smaller pseudo-count. All other components remain unchanged: using the same internal representations, we select the action whose representation exhibits the smallest novelty as the final output. Table 7 reports the success rate comparison. As shown, using CFN as the pseudo-counter
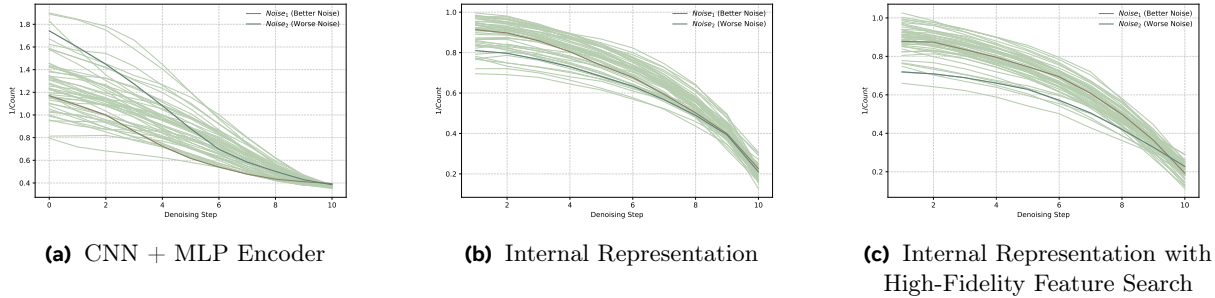
**(a)** CNN + MLP Encoder    **(b)** Internal Representation    **(c)** Internal Representation with High-Fidelity Feature Search

**Figure 8** Comparison of the Count evolution during the action denoising process under different representation extraction methods, illustrated with the Block Handover task as an example.

**Table 7** Performance comparison of final success rates (%) using different pseudo-counters

|  | Block Handover | Container Place | Diverse Bottles Pick | Average |
|---|---|---|---|---|
| $\pi_0$ | 41.0 | 25.0 | 21.0 | 29.0 |
| $\pi_0$ **+ TACO (RND)** | 54.0 | 33.0 | **30.0** | 39.0 |
| $\pi_0$ **+ TACO (CFN, Ours)** | **62.0** | **40.0** | 27.0 | **43.0** |

yields consistently better average performance.

# G    Additional Real-World Experiment Details

**Real-World Platform and Task Setup.** As shown in Figure 9, our platform consists of a Realman RM75-6F dual-arm robot, two Robotiq 2F grippers, one Intel RealSense L515 camera serving as the main-view camera, and two Intel RealSense D405 cameras mounted on the wrists. including one main-view camera and two wrist-mounted cameras. A workstation equipped with an RTX 4090 GPU is used for model deployment and inference.

The prompts for each task are listed in Table 8. During testing, each object is randomly placed within an area of approximately $3\,\text{cm} \times 4\,\text{cm}$. The overall task execution process is illustrated in Figure 9.
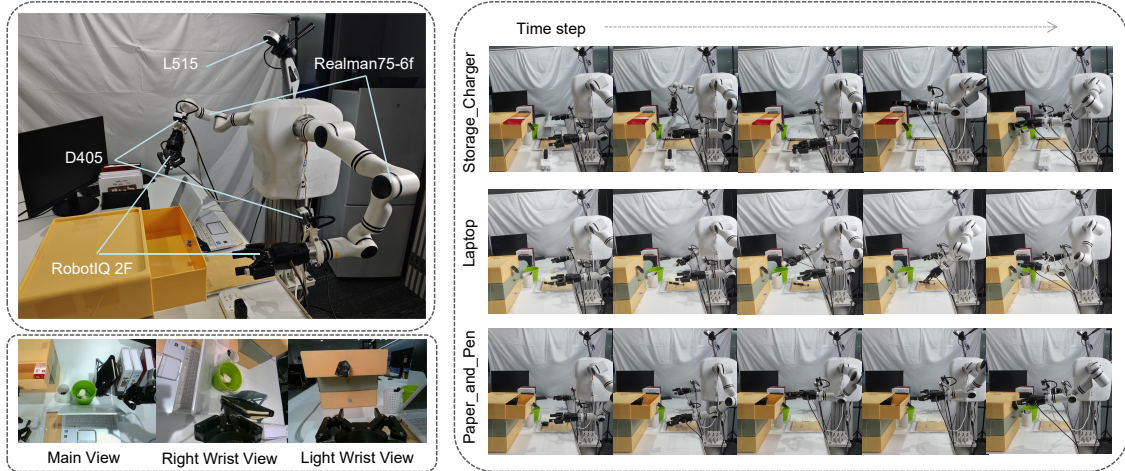
**Figure 9** Schematic illustration of our real-world robotic platform (left), and the task execution pipeline for the three real-world tasks (right).

| Task | Prompt |
|---|---|
| Receive Book | Use the right arm to catch the book handed over by the person. |
| Paper and Pen | Place the paper ball into the trash bin with the left arm, and put the pen into the pen holder with the right arm. |
| Storage Charger | The left arm places the charging head into the drawer and closes it. |
| Pick Books | Grip the book with both hands, one hand on each side, and lift it simultaneously with both arms. |
| Laptop | Close the laptop lid using the right arm, unplug the charger with the left arm, place it in the top drawer, and close the drawer. |

**Table 8** Tasks in real world and their corresponding prompts.