
Efficient Reinforcement Learning by Guiding Generalist World Models with Non-Curated Data

Yi Zhao[†] ¹ **Aidan Scannell**^{1,2} **Wenshuai Zhao**¹ **Yuxin Hou**³ **Tianyu Cui**^{1,4}

Le Chen⁵ **Dieter Büchler**^{5,6,7} **Arno Solin**¹ **Juho Kannala**^{1,8} **Joni Pajarinen**¹

¹Aalto University ²University of Edinburgh ³Deep Render ⁴Imperial College London

⁵Max Planck Institute for Intelligent Systems ⁶University of Alberta

⁷Alberta Machine Intelligence Institute (Amii) ⁸University of Oulu

Abstract

Leveraging offline data is a promising way to improve the sample efficiency of online reinforcement learning (RL). This paper expands the pool of usable data for offline-to-online RL by leveraging abundant non-curated data that is reward-free, of mixed quality, and collected across multiple embodiments. Although learning a world model appears promising for utilizing such data, we find that naive fine-tuning fails to accelerate RL training on many tasks. Through careful investigation, we attribute this failure to the distributional shift between offline and online data during fine-tuning. To address this issue and effectively use the offline data, we propose two essential techniques: *i*) experience rehearsal and *ii*) execution guidance. With these modifications, the non-curated offline data substantially improves RL’s sample efficiency. Under limited sample budgets, our method achieves a 102.8% relative improvement in aggregate score over learning-from-scratch baselines across 72 visuomotor tasks spanning 6 embodiments. On challenging tasks such as locomotion and robotic manipulation, it outperforms prior methods that utilize offline data by a decent margin.

1 Introduction

Leveraging offline data offers a promising way to improve the sample efficiency of Reinforcement Learning (RL). Prior work has focused primarily on utilizing curated offline data labeled with rewards [1–4], which is expensive and laborious to obtain. For instance, leveraging offline datasets for new robotic manipulation tasks requires retrospectively annotating image-based data with rewards. We instead propose expanding the pool of usable offline data by utilizing abundant non-curated data that is reward-free, of mixed quality, and collected across multiple embodiments. This leads to our primary research question:

How can we effectively leverage non-curated offline data for efficient RL?

Typical offline-to-online RL [5–9] methods fail to utilize non-curated offline data due to their assumption of structured data with rewards. While pre-training visual encoders [10–15] is a common approach to utilize non-curated offline datasets, it fails to fully leverage the rich information, such as dynamics models, informative states, and action priors. On the other hand, learning world models from offline data appears promising for utilizing the non-curated dataset. However, prior work has explored world model training primarily in settings with known rewards [16–18] or expert demonstrations [19–21] or focused solely on visual prediction [22, 23]. Recent approaches [24–26] have developed novel architectures for world model pre-training using in-the-wild action-free data, but paid limited attention to the fine-tuning process. As a result, despite being trained on massive datasets, these methods show

[†] Correspondence to yi.zhao@aalto.fi.

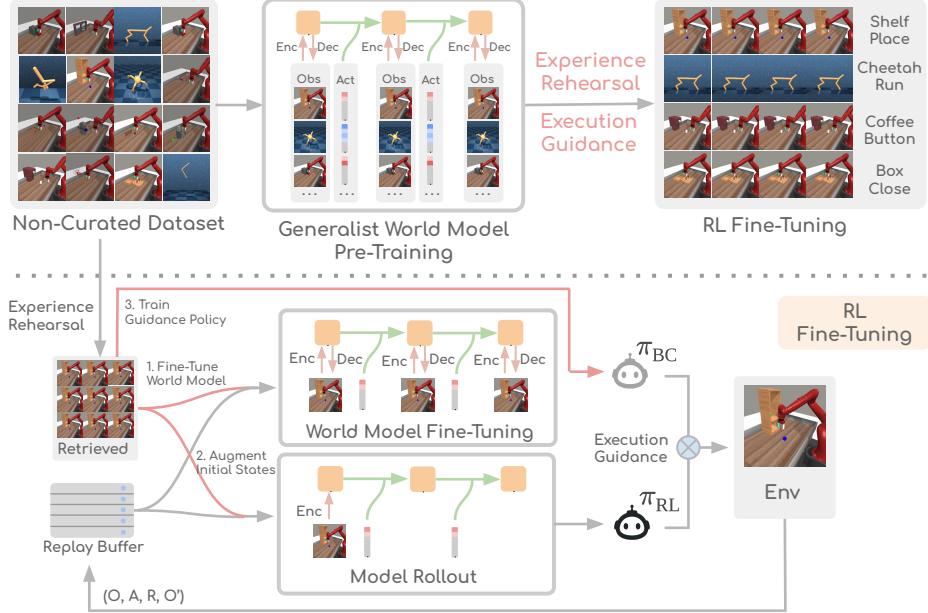


Figure 1: Overview of GSA (Generalist-to-Specialist Adaptation). Facing non-curated offline data with reward-free, mixed-quality, and multi-embodiment data, we train a task and embodiment-agnostic world model. Combined with experience rehearsal and execution guidance, the pre-trained world model improves the sample efficiency of RL training over a wide range of tasks.

only marginal improvements over training-from-scratch baselines. Additionally, due to the computational costs of RL experiments, previous work evaluated only on a small set of tasks, leaving the effectiveness of the learned world model unclear on broader tasks. In contrast, we extensively evaluate our method on 72 visuomotor control tasks – 12 times more than previous baselines – spanning both locomotion and robotic manipulation, demonstrating consistent improvements over existing approaches.

Through experiments, we observe that naively fine-tuning a world model fails to improve RL’s sample efficiency on many tasks. Through careful investigation, we identify the root cause as a distributional shift between offline and online data during fine-tuning. Specifically, when the offline data distribution does not sufficiently cover the data distribution of downstream tasks, the pre-trained world models struggle to benefit policy learning due to this distribution mismatch. Especially for tasks requiring hard exploration, even when the pre-trained world model performs well on “promising” state-action distributions, the agent often fails to learn because it cannot effectively reach these promising regions. Building on these insights, we propose using offline data in both pre-training and fine-tuning stages, in contrast to previous methods that only consider the offline data for world model pre-training [26–28]. By this end, we propose Generalist-to-Specialist Adaptation (GSA), which enables generalist world models pre-trained on non-curated data to master a wide range of downstream tasks. GSA has two essential techniques: experience rehearsal and execution guidance. Experience rehearsal mitigates distributional shift by retrieving task-relevant trajectories from offline datasets, while execution guidance promotes exploration by steering the agent toward regions where the world model has high confidence.

Equipped with our proposed techniques, GSA demonstrates strong performance across a diverse set of tasks. Specifically, under a limited sample budget (150k samples), GSA achieves a 102.8% relative improvement in aggregate score over learning-from-scratch baselines (DrQ v2 and Dreamer v3), while matching their performance achieved with larger sample budgets. On representative challenging tasks, GSA outperforms baselines that leverage offline data as well as state-of-the-art methods using pre-trained world models by a significant margin. Additionally, without any modifications, we show that GSA improves task adaptation, enabling agents to efficiently adapt their skills to new tasks.

To summarize, our contributions are:

- (C1) We propose a more realistic setting for leveraging offline data that consists of reward-free and mixed-quality multi-embodiment data.

Table 1: Comparison with different policy learning methods that leverage offline data.

	Offline RL	Off2On RL	RLPD	MT Offline RL	GSA (ours)
Reward-free offline data	✗	✗	✗	✗	✓
Non-expert offline data	✓	✓	✓	✓	✓
X-embodiment offline data	✗	✗	✗	✓	✓
Continual Improvement	✗	✓	✓	✗	✓
Training stability	✗	✗	✗	✗	✓

- (C2) We demonstrate that naive world model fine-tuning fails on many tasks due to distributional shift between pre-training and fine-tuning data.
- (C3) We propose two essential techniques, experience rehearsal and execution guidance, to mitigate the distributional gap and encourage exploration.
- (C4) We present GSA, which effectively leverages non-curated offline data and significantly outperforms existing approaches across a diverse set of tasks.

2 Related Work

In this section, we review methods that leverage offline data, including pre-training in the context of RL and world models. See App. C for more related work and Table 1 for an overview comparison.

RL with task-specific offline datasets Leveraging offline data is a promising direction to improve sample efficiency in RL. One representative approach is offline RL, which trains agents using offline data without environment interaction. These methods typically constrain the distance between learned and behavior policies in different ways [2, 3, 29–33]. However, policy performance is highly dependent on dataset quality [34]. To enable continued improvement, offline-to-online RL methods [5–7, 9, 17] were developed, which fine-tune policies trained with offline RL by interacting with environments. MOTO [17] proposes a model-based offline-to-online RL method with reward-labeled data, and requires model-based value expansion, policy regularization, and controlling epistemic uncertainty, while our method leverages reward-free and multi-embodiment data and requires none of the techniques proposed by MOTO.

Typical offline-to-online RL face training instability challenges [5, 16]. To mitigate this issue, RLPD [35] is proposed and demonstrates strong performance by simply concatenating offline and online data, but requires reward-labeled task-specific offline data and does not address multi-embodiment scenarios. ExPLORe [36] labels reward-free offline data using approximated upper confidence bounds (UCB) to solve hard exploration tasks, but relies on near-expert data for the target tasks, while we consider a more general setting with non-curated data.

RL with multi-task offline datasets Recent work has explored multi-task offline RL [4, 18, 37–39], but requires known rewards. PWM [40] and TDMPC v2 [18] train world models for multi-task RL but are limited to state-based inputs and reward-labeled data. To handle unknown rewards, approaches like human labeling [41, 42], inverse RL [43, 44], or generative adversarial imitation learning [45] can be used, though these require human labor or expert demonstrations. Yu et al. [46] assigns zero rewards to unlabeled data, which introduces additional bias. Apart from these, there is a line of work that focuses on representation learning from in-the-wild data [10, 12, 14, 27, 47–52] but fails to utilize rich information in the dataset, such as dynamics.

Recent studies [24, 26, 53] explore world model pre-training with action-free data, focusing on world model architecture design to utilize the action-free data. However, we demonstrate that naive fine-tuning of pre-trained world models fails on challenging tasks, while our method, incorporating experience rehearsal and execution guidance, significantly improves RL performance across 72 tasks.

3 Methods

In this section, we detail our two-stage approach, which consists of (*i*) world model pre-training, which learns a multi-task & embodiment world model, given offline data, which rather importantly,

includes reward-free, mixed-quality data, and (ii) RL-based fine-tuning which leverages the pre-trained world model and online interaction in an offline-to-online fashion. See Fig. 1 for the overview and Alg. 1 for the full algorithm.

3.1 Problem Setup

In this paper, we assume the agent has access to a non-curated but in-domain offline dataset \mathcal{D}_{off} with three key characteristics: (i) trajectories lack reward labels r_t^i , (ii) data quality is mixed, and (iii) data comes from multiple embodiments. During fine-tuning, the agent interacts with the environment to collect labeled trajectories $\tau_{\text{on}}^i = \{o_t^i, a_t^i, r_t^i\}_{t=1}^T$ and stores them in an online dataset $\mathcal{D}_{\text{on}} = \{\tau_{\text{on}}^i\}_{i=1}^{N_{\text{on}}}$. Our goal is to learn a high-performance policy by leveraging both \mathcal{D}_{off} and \mathcal{D}_{on} while minimizing the required online interactions N_{on} .

3.2 Multi-Embodiment World Model Pre-training

During pre-training, rather than training separate models per task as in previous work [54–56], we train one world model per benchmark and demonstrate that a single multi-task & embodiment world model can effectively leverage non-curated data.

Since our primary goal is enabling RL agents to use non-curated offline data rather than proposing a new architecture, we adopt the widely-used recurrent state space model (RSSM) [57] with several modifications: (i) removal of task-related losses, (ii) zero-padding of actions to unify dimensions across embodiments, and (iii) scaling the model to 280M parameters. With these changes, we show that RSSMs can successfully learn the dynamics of multiple embodiments and can be fine-tuned for various tasks.

Our first stage pre-trains the following components:

$$\begin{array}{ll} \text{Sequence model : } h_t = f_\theta(h_{t-1}, z_{t-1}, a_{t-1}) & \text{Encoder : } z_t \sim q_\theta(z_t | h_t, o_t) \\ \text{Dynamics predictor : } \hat{z}_t \sim p_\theta(z_t | h_t) & \text{Decoder : } \hat{o}_t \sim d_\theta(\hat{o}_t | h_t, z_t). \end{array}$$

The models f_θ , q_θ , p_θ and d_θ are jointly optimized by minimizing:

$$\mathcal{L}(\theta) = \mathbb{E}_{p_\theta, q_\theta} \left[\sum_{t=1}^T \underbrace{-\ln p_\theta(o_t | z_t, h_t)}_{\text{pixel reconstruction loss}} + \beta \cdot \underbrace{\text{KL}(q_\theta(z_t | h_t, o_t) \| p_\theta(z_t | h_t))}_{\text{latent state consistency loss}} \right]. \quad (1)$$

The first term minimizes reconstruction error while the second enables latent dynamics learning. While there is room to improve world model pre-training through recent self-supervised methods [58] or advanced architectures [59, 60], such improvements are orthogonal to our method and left for future work.

3.3 RL-based Fine-Tuning with Rehearsal and Guidance

In our fine-tuning stage, the agent interacts with the environment to collect new data $\tau_{\text{on}}^i = \{o_t^i, a_t^i, r_t^i\}_{t=0}^T$. This data is used to learn a reward function via supervised learning while fine-tuning the world model with Eq. (1). For simplicity, we represent the concatenation of h_t and z_t as $s_t = [h_t, z_t]$. The actor and critic are trained using imagined trajectories $\tilde{\tau}$ generated by rolling out the policy $\pi_\phi(a | s)$ in the world model p_θ , starting from initial states $p_0(s)$ sampled from the replay buffer. The critic $v_\phi(V_t^\lambda | s_t)$ learns to approximate the distribution over the λ -return V_t^λ , calculated as:

$$\underbrace{V_t^\lambda}_{\lambda-\text{return}} = \hat{r}_t + \gamma \begin{cases} (1 - \lambda)v_{t+1}^\lambda + \lambda V_{t+1}^\lambda & \text{if } t < H \\ v_H^\lambda & \text{if } t = H \end{cases} \quad (2)$$

where $v_t^\lambda = \mathbb{E}[v_\phi(V_t^\lambda | s_t)]$ denotes the critic's expected value predicted. The value function v_ϕ and actor π_ϕ are updated by maximizing the log likelihood and entropy-regularized λ -return respectively:

$$\mathcal{L}(v_\phi) = \mathbb{E}_{p_\theta, \pi_\phi} \left[- \sum_{t=1}^{H-1} \ln v_\phi(V_t^\lambda | s_t) \right], \quad \mathcal{L}(\pi_\phi) = \mathbb{E}_{p_\theta, \pi_\phi} \left[\sum_{t=1}^{H-1} (-v_t^\lambda - \eta \cdot \mathbf{H}[a_t | s_t]) \right]. \quad (3)$$

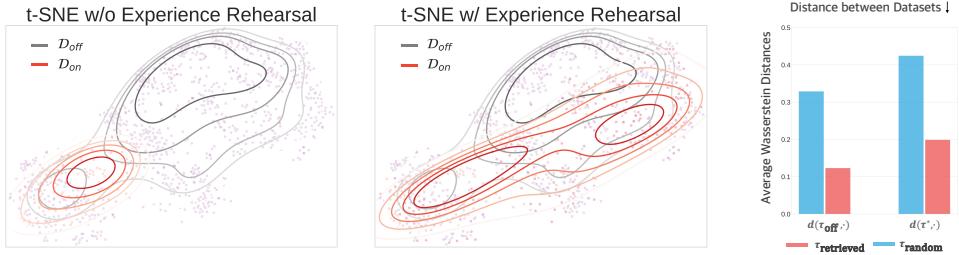


Figure 2: **Visualization of Distribution Mismatch.** **Left:** At the early stage of fine-tuning, there is a distribution shift between offline data used for world model pre-training and online data used for RL fine-tuning, which hurts performance. **Middle:** Experience rehearsal mitigates the distributional shift issue. **Right:** Quantitatively, at the early stage of fine-tuning, experience rehearsal reduces the Wasserstein distance between the online data and both the offline and expert data.

Motivation While previous methods typically discard non-curated offline data during fine-tuning [26, 28, 53], we find that relying solely on a pre-trained world model often fails, particularly on hard-exploration tasks. To understand why, we analyze the Shelf Place task from Meta-World as an illustrative task by visualizing the distributions of offline data \mathcal{D}_{off} used for world model pre-training and online data \mathcal{D}_{on} collected during early RL training in Fig. 2. The t-SNE plot in Fig. 2 (left) reveals a distribution mismatch between \mathcal{D}_{off} and \mathcal{D}_{on} , leading to three key issues: (i) The world model’s accuracy degrades on states visited by the early-stage policy, especially when the offline data distribution is narrow, hurting sample efficiency. (ii) For hard exploration tasks, the agent struggles to reach high-reward regions, causing the world model to be fine-tuned on a narrow online data distribution and leading to catastrophic forgetting. (iii) The policy update in Eq. (3) relies on imagined trajectories $\tilde{\tau} = p_0(s) \prod_{t=0}^{H-1} \pi_\phi(a_t | s_t) p_\theta(s_{t+1} | s_t, a_t)$, where $p_0(s)$ is sampled from \mathcal{D}_{on} . A narrow $p_0(s)$ limits the world model to rollout promising trajectories for policy updates. To address these challenges, we introduce two key components: *i*) experience rehearsal, which mitigates distributional shift by retrieving task-relevant trajectories from non-curated datasets (Fig. 2 middle, right), and *ii*) execution guidance, which encourages exploration by steering the agent toward regions where the world model has high confidence.

Experience Rehearsal Prior work like RLPD [35] and ExPLORe [36] has shown that replaying offline data can boost RL training. However, these methods use small, well-structured offline datasets. In our setting, directly replaying non-curated offline data is infeasible since our datasets are $\sim 100\times$ larger and contain diverse tasks and embodiments.

We propose retrieving task-relevant trajectories $\mathcal{D}_{\text{retrieved}} = \{\tau_{\text{retrieved}}^i\}_{i=1}^N$ from the non-curated offline data based on neural feature distance between online samples and offline trajectories. This filters out irrelevant trajectories, creating a small dataset of task-relevant data. Specifically, we compute:

$$\mathbf{D} = \|\mathbf{e}_\theta(o_{\text{on}}) - \mathbf{e}_\theta(o_{\text{off}})\|_2, \quad (4)$$

where \mathbf{e}_θ is the encoder learned during world model pre-training, and o_{on} and o_{off} are initial observations from trajectories in the online buffer and offline dataset, respectively. For efficient search the top-k similar trajectories, we pre-compute key-value pairs mapping trajectory IDs to neural features and use Faiss [61], enabling retrieval in seconds.

The retrieved data is replayed during fine-tuning, so-called experience rehearsal. The retrieved data serves three purposes, as shown in Fig. 1. First, it prevents catastrophic forgetting by continuing to train the world model on relevant pre-training data, particularly important for hard exploration tasks with narrow online data distributions. Second, it augments the initial state distribution $p_0(s)$ during model rollout, enabling the world model to rollout promising trajectories for policy learning. Third, as described below, it enables learning a policy prior for execution guidance. Unlike RLPD and ExPLORe, we do not use this data to learn a Q-function, eliminating the need for reward labels.

Execution Guidance via Prior Actors Standard RL training initializes the replay buffer with random actions and collects new data through environment interaction using the training policy. However, offline data often contains valuable information like near-expert trajectories and diverse state-action coverage that should be utilized during fine-tuning. Additionally, distribution shift

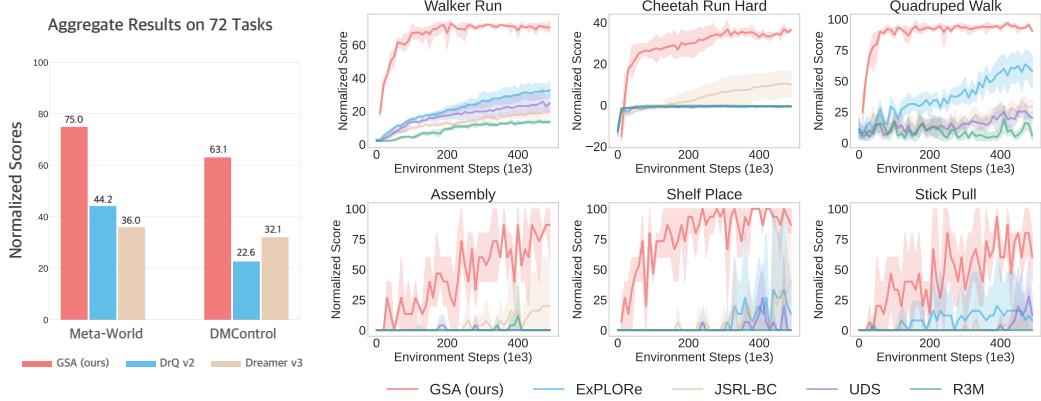


Figure 3: **Left:** Quantitative comparison across 72 diverse tasks from Meta-World and DMControl. GSA achieves a 102.8% relative improvement in aggregate score over learning-from-scratch baselines when using the same sample budget (150k). It also matches the performance of baselines even when they are trained with substantially more samples (see App. I for full results). **Right:** Learning curves on representative challenging locomotion and robotic manipulation tasks. GSA consistently outperforms state-of-the-art methods that leverage offline data by a decent margin. We plot the mean and corresponding 95% confidence interval.

between offline and online data can degrade pre-trained model weights, making it important to guide the online data collection toward the offline distribution at the early training stage.

To achieve this, we train a prior policy π_{bc} via behavioral cloning on the retrieved offline data $\mathcal{D}_{\text{retrieved}}$. During online data collection, we alternate between this prior policy π_{bc} and the RL policy π_ϕ according to a pre-defined schedule. Specifically, at the start of each episode, we probabilistically select whether to use π_{bc} . If π_{bc} is selected, we randomly choose a starting timestep t_{bc} and duration H during which π_{bc} is active, with π_ϕ used for the rest timesteps.

While this approach shares similarities with JSRL [33], our method differs in three key aspects: *i*) we leverage non-curated rather than task-specific offline data, *ii*) we demonstrate the benefits of a model-based approach over JSRL’s model-free framework, and *iii*) we randomly switch between policies mid-episode rather than only using π_{bc} at episode start. The complete algorithm and theoretical analysis can be found in App. H and App. B, respectively.

4 Experiments

In the experiments, we aim to answer the following questions: *(i)* How does GSA compare to state-of-the-art methods that leverage offline data and train-from-scratch baselines in terms of sample efficiency and final performance? *(ii)* How does GSA compare to other leading model-based approaches that utilize offline data? *(iii)* How effectively does GSA adapt to new tasks in a continual learning setting? We further conduct detailed ablation studies to thoroughly evaluate our method.

Tasks We evaluate our method on *pixel*-based continuous control tasks from DMControl and Meta-World. The chosen tasks include both locomotion and manipulation tasks covering different challenges in RL, including high-dimensional observations, hard exploration, and complex dynamics. We use three random seeds for each task.

Dataset Our dataset consists of data from two benchmarks: DeepMind Control Suite (DMControl) and Meta-World, visualized in App. K. For DMControl, we include 10k trajectories covering 5 embodiments collected by *unsupervised RL agents* [28, 62], trained via curiosity without task-related information. For Meta-World, we collect 50k trajectories across 50 tasks using pre-trained RL agents from TDMPC v2 [18]. To create a diverse dataset with varying quality, we inject Gaussian noise with $\sigma \in \{0.1, 0.3, 0.5, 1.0, 2.0\}$ during execution. Combined with the DMControl data, our complete offline dataset contains 60k trajectories (10M state-action pairs) spanning 6 embodiments.

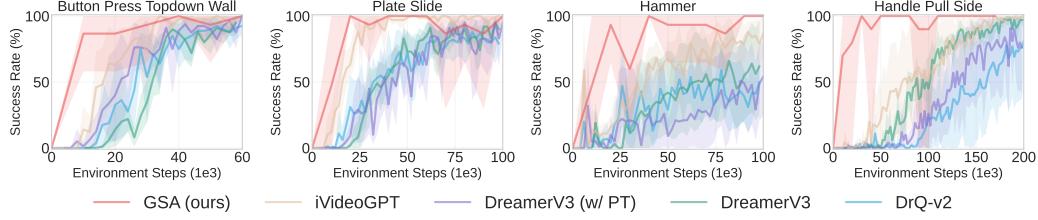


Figure 4: Comparison with other world model pre-training methods. GSA outperforms state-of-the-art model-based methods without relying on techniques used in iVideoGPT, such as reward shaping and demonstration-based replay buffer initialization.

4.1 GSA Improves Sample Efficiency Across Diverse Tasks

Comparison with Methods that Leverage Offline Data We compare GSA against several state-of-the-art methods that leverage reward-free data to improve RL training: (i) **R3M** [11], a visual representation pre-training approach that serves as our baseline for comparing pre-trained visual features using non-curated offline data. (ii) **UDS** [46], which assigns zero rewards to offline data. Since UDS is an offline RL method, we combine it with RLPD [35] for policy training. (iii) **ExPLORe** [36], which labels offline data using UCB rewards. We enhance the original implementation with reward ensembles. (iv) **JSRL-BC** [33], which collects online data using a mixture of the training policy and a behavior-cloned prior policy learned from offline data. As the compared baselines cannot handle multi-embodiment data like GSA, we preprocess the offline data to only include task-relevant trajectories for them. Despite the baselines having access to better-structured data, GSA still significantly outperforms all baselines across the tested tasks. See App. G for the details of baselines.

Results Fig. 3 (right) shows comparison results with baselines. Our method outperforms *all* compared baselines by a large margin. Compared to R3M, GSA shows the importance of world model pre-training and reusing offline data during fine-tuning, versus representation learning alone. R3M fails to improve sample efficiency on most tasks, consistent with findings in Hansen et al. [63].

UDS and ExPLORe reuse offline data by labeling it with zero rewards and UCB rewards, respectively, and concatenating it with online data for off-policy updates. UDS shows only slightly better performance on Walker Run compared to R3M and JSRL-BC, demonstrating the ineffectiveness of zero-reward labeling. ExPLORe performs better on 2/3 locomotion tasks and shows progress on challenging manipulation tasks, but GSA still significantly outperforms it, demonstrating the superiority of leveraging a generalist world model and properly reusing offline data during fine-tuning.

GSA also clearly outperforms JSRL-BC. JSRL-BC’s performance heavily depends on offline data distribution. While JSRL-BC can perform well when a good prior actor can be extracted from offline data, it struggles with non-expert trajectories, showing only marginal improvements over other baselines on Cheetah Run Hard, Assembly, and Shelf Place tasks. In contrast, GSA effectively leverages non-expert offline data. For example, on Quadruped Walk, GSA benefits from exploratory offline data collected by unsupervised RL, enabling pixel-based control within just 100 trials.

Comparison with Training-from-Scratch Methods We compare GSA with two widely used training-from-scratch baselines: **DrQ v2** and **Dreamer v3**, representing model-free and model-based approaches, respectively. Figure 3 (left) and App. I show comparison results on 22 locomotion and 50 robotic manipulation tasks with pixel inputs from DMControl and Meta-World benchmarks. With 150k online samples, GSA achieves a 102.8% relative improvement in aggregate score compared to DrQ v2 and Dreamer v3, matching their performance obtained with 3.3-6.7 \times more samples (500k for DMControl, 1M for Meta-World). For example, GSA enables an Ant robot to walk forward within 100 trials, while widely used learning-from-scratch baselines require **10-30 \times** more samples. Furthermore, GSA achieves promising performance on hard exploration tasks where learning-from-scratch baselines fail, such as challenging Meta-World manipulation tasks and hard DMControl tasks.

Comparison with Other Model-Based Methods While most multi-task/multi-embodiment world models focus on visual prediction [22, 23] or imitation learning [19, 20], some works like Seo et al. [24], Wu et al. [25], and iVideoGPT [26] investigate world model pre-training with in-the-wild data for RL. These methods typically focus on designing novel or scalable model architectures to leverage the offline data, but lack mechanisms to better leverage offline data during RL fine-tuning.

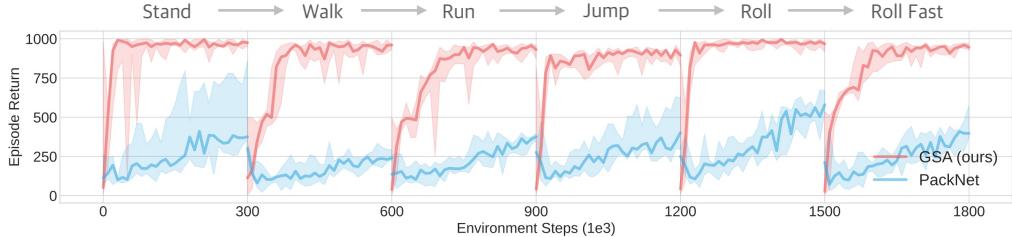


Figure 5: GSA enables fast task adaptation. We train an RL agent to control an Ant robot from DMControl to complete a series of tasks incrementally. GSA significantly outperforms the widely used baseline PackNet by properly leveraging non-curated offline data.

Furthermore, due to the cost of RL training, these methods are usually evaluated on limited task sets, making the effectiveness of the pre-trained world model unclear on diverse tasks.

Figure 4 compares our method with world model pre-training approaches. The baseline results are from the iVideoGPT paper. Despite extensive pre-training on diverse manipulation data, iVideoGPT and pre-trained DreamerV3 show only marginal improvements over training-from-scratch baselines. In contrast, GSA clearly accelerates RL training by properly leveraging non-curated offline data during both pre-training and fine-tuning. Notably, baselines in Fig. 4 use reward shaping and expert replay buffer pre-filling, while GSA uses *none* of these tricks yet achieves superior performance. This highlights that *(i)* non-curated offline data contains useful information for RL fine-tuning, and *(ii)* GSA can effectively leverage such data. Furthermore, GSA could potentially be combined with iVideoGPT to leverage even more diverse offline data in future work.

4.2 GSA Enables Fast Task Adaptation

We investigate GSA’s benefits for continual task adaptation, where an agent must incrementally solve a sequence of tasks. While similar to continual reinforcement learning (CRL) or life-long RL [64, 65], we use a simplified setting with a limited task set. Note that CRL has a broad scope; assumptions and experiment setups vary among methods, making it difficult to set up a fair comparison with other methods. Rather than proposing a state-of-the-art CRL method, we aim to demonstrate that GSA offers an effective approach to leverage previous data that also fits the CRL setting.

Setup & Baselines We set our continual multi-task adaptation experiment based on the Ant robot from the DMControl. Specifically, the agent sequentially learns stand, walk, run, jump, roll, and roll fast tasks with 300K environment steps per task. To have a fair comparison, i.e., having comparable model parameters and eliminating the potential effects from pre-training on other tasks, we pre-train a small world model only on the Ant domain. During training, the agent can access all previous experiences and model weights. We compare against a widely used baseline PackNet [66], which iteratively prunes actor parameters while preserving important weights (with larger magnitude) for previous skills. For each new task, PackNet fine-tunes the actor model via iterative pruning while randomly reinitializing the critic model since rewards are not shared among tasks.

Results Figure 5 shows GSA significantly outperforms PackNet, enabling adaptation within 100 trials per task. With limited samples, PackNet achieves only 20–60% of GSA’s episodic returns. We attribute GSA’s superior performance to its ability to leverage the diverse offline data through both world model pre-training and fine-tuning with experience rehearsal, and execution guidance.

4.3 Ablations

Role of Each Component We now analyze each component’s contribution using the same set of tasks from Sec. 4.1. As shown in Fig. 6, world model pre-training shows promising results when the offline data consists of diverse trajectories, such as data collected by exploratory agents (Walker Run), while it fails to work well when the offline data distribution is relatively narrow as in the Meta-World tasks. We found that experience rehearsal and execution guidance stabilize training and improve performance on hard exploration tasks like Cheetah Run Hard and challenging manipulation tasks from Meta-World. This addresses *(i)* world model pre-training alone, failing to fully leverage rich state and action information from the non-curated offline data and *(ii)* distributional shift between

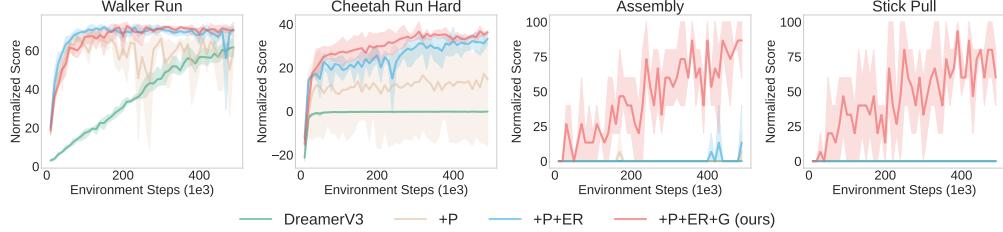


Figure 6: Ablation study on key components. “P” represents world model pre-training, “ER” means experience rehearsal, and “G” represents execution guidance. The combination of pre-trained generalist world models with retrieval-based experience rehearsal and execution guidance boosts RL performance across diverse tasks.

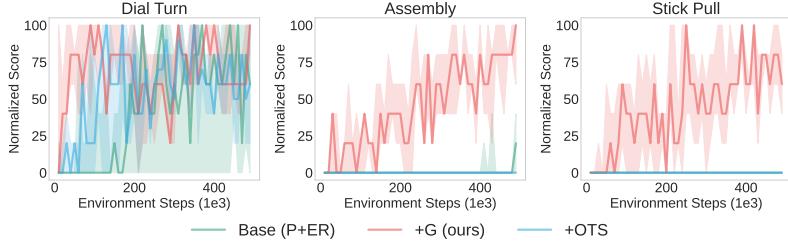


Figure 7: Comparison of execution guidance versus uncertainty-based reward labeling. GSA demonstrates the effectiveness of using execution guidance over uncertainty-based reward labeling on challenging robotic manipulation tasks.

offline and online data during fine-tuning hurts the learning. The proposed retrieval-based experience rehearsal and execution guidance help utilize offline data and accelerate exploration, which together enable GSA to achieve strong performance on a wide range of tasks.

Comparison with Uncertainty-Aware Reward Function To leverage reward-free offline data, ExpLORE [36] proposes to label offline data with uncertainty-based rewards. To demonstrate the effectiveness of GSA, we compare it with uncertainty-based rewards. Specifically, instead of using execution guidance, we use Optimistic Thompson Sampling (OTS) [67] to label the imagined trajectories via model rollout. As shown in Fig. 7, our method outperforms the variant using OTS on hard exploration tasks, Assembly and Stick Pull, by a large margin, showing the effectiveness of using execution guidance.

Comparison of Fine-Tuning Different Components We now investigate the role of different components in the world model during fine-tuning. We use the Quadruped Walk task as a representative task for the investigation. As shown in Fig. 8, encoder and decoder as well as latent dynamics play compatible roles during fine-tuning. Fine-tuning the full world model yields the best performance on the tested task. The full world model is fine-tuned by default in our experiments.

5 Conclusion

We propose GSA, a simple yet efficient approach to leverage ample non-curated offline datasets consisting of reward-free, mixed-quality data collected across multiple embodiments. GSA pre-trains a generalist world model on the non-curated data and adapts to downstream tasks via RL. We show that naive fine-tuning of world models fails to accelerate RL training due to distributional shift and propose two essential techniques - experience rehearsal and execution guidance - to mitigate this issue. Equipped with these techniques, we demonstrate that generalist world models pre-trained on non-curated data are able to boost RL’s sample efficiency across a broader range of locomotion and robotic manipulation tasks. We compared GSA against a wide set of baselines, including two widely used training-from-scratch methods, five methods that utilize offline data, and one continual

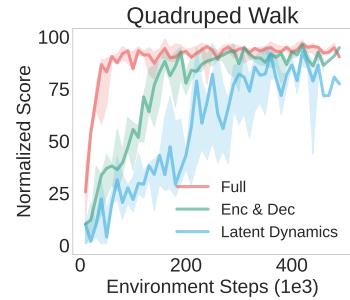


Figure 8: Impact of fine-tuning different world model components.

learning method. Our GSA consistently delivers strong performance over these baselines. Extensive ablation studies reveal the effectiveness of the proposed techniques. While promising, GSA can be improved in multiple ways: extending to real-world applications, leveraging in-the-wild offline data, and exploring novel world model architectures.

Acknowledgments

We acknowledge CSC – IT Center for Science, Finland, for awarding this project access to the LUMI supercomputer, owned by the EuroHPC Joint Undertaking, hosted by CSC (Finland) and the LUMI consortium through CSC. We acknowledge the computational resources provided by the Aalto Science-IT project. We acknowledge funding from the Research Council of Finland (353138, 362407, 352788, 357301, 339730). Aidan Scannell was supported by the Research Council of Finland, Flagship program Finnish Center for Artificial Intelligence (FCAI).

References

- [1] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline Reinforcement Learning: Tutorial, Review, and Perspectives on Open Problems. *arXiv preprint arXiv:2005.01643*, 2020.
- [2] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative Q-Learning for Offline Reinforcement Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [3] Scott Fujimoto and Shixiang Shane Gu. A Minimalist Approach to Offline Reinforcement Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [4] Aviral Kumar, Anikait Singh, Frederik Ebert, Mitsuhiro Nakamoto, Yanlai Yang, Chelsea Finn, and Sergey Levine. Pre-Training for Robots: Offline RL Enables Learning New Tasks from a Handful of Trials. In *Robotics: Science and Systems (RSS)*, 2023.
- [5] Seunghyun Lee, Younggyo Seo, Kimin Lee, Pieter Abbeel, and Jinwoo Shin. Offline-to-Online Reinforcement Learning via Balanced Replay and Pessimistic Q-Ensemble. In *Conference on Robot Learning (CoRL)*, 2022.
- [6] Yi Zhao, Rinu Boney, Alexander Ilin, Juho Kannala, and Joni Pajarinen. Adaptive Behavior Cloning Regularization for Stable Offline-to-Online Reinforcement Learning. *arXiv preprint arXiv:2210.13846*, 2022.
- [7] Zishun Yu and Xinhua Zhang. Actor-Critic Alignment for Offline-to-Online Reinforcement Learning. In *International Conference on Machine Learning (ICML)*, 2023.
- [8] Mitsuhiro Nakamoto, Simon Zhai, Anikait Singh, Max Sobol Mark, Yi Ma, Chelsea Finn, Aviral Kumar, and Sergey Levine. Cal-QL: Calibrated Offline RL Pre-Training for Efficient Online Fine-Tuning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- [9] Ashvin Nair, Abhishek Gupta, Murtaza Dalal, and Sergey Levine. AWAC: Accelerating Online Reinforcement Learning with Offline Datasets. *arXiv preprint arXiv:2006.09359*, 2020.
- [10] Max Schwarzer, Nitarshan Rajkumar, Michael Noukhovitch, Ankesh Anand, Laurent Charlin, R Devon Hjelm, Philip Bachman, and Aaron C Courville. Pretraining Representations for Data-Efficient Reinforcement Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [11] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3M: A Universal Visual Representation for Robot Manipulation. In *Conference on Robot Learning (CoRL)*, 2022.
- [12] Simone Parisi, Aravind Rajeswaran, Senthil Purushwalkam, and Abhinav Gupta. The Unsurprising Effectiveness of Pre-Trained Vision Models for Control. In *International Conference on Machine Learning (ICML)*, 2022.
- [13] Tete Xiao, Ilija Radosavovic, Trevor Darrell, and Jitendra Malik. Masked Visual Pre-Training for Motor Control. *arXiv preprint arXiv:2203.06173*, 2022.
- [14] Mengjiao Yang and Ofir Nachum. Representation Matters: Offline Pretraining for Sequential Decision Making. In *International Conference on Machine Learning (ICML)*, 2021.
- [15] Jinghuan Shang, Karl Schmeckpeper, Brandon B May, Maria Vittoria Minniti, Tarik Kelestemur, David Watkins, and Laura Herlant. Theia: Distilling Diverse Vision Foundation Models for Robot Learning. In *Conference on Robot Learning (CoRL)*, 2024.
- [16] Cong Lu, Philip J Ball, Tim GJ Rudner, Jack Parker-Holder, Michael A Osborne, and Yee Whye Teh. Challenges and Opportunities in Offline Reinforcement Learning from Visual Observations. *Transactions*

on Machine Learning Research (TMLR), 2023.

- [17] Rafael Rafailov, Kyle Beltran Hatch, Victor Kolev, John D Martin, Mariano Phiellipp, and Chelsea Finn. MOTO: Offline Pre-Training to Online Fine-tuning for Model-Based Robot Learning. In *Conference on Robot Learning (CoRL)*, 2023.
- [18] Nicklas Hansen, Hao Su, and Xiaolong Wang. TD-MPC2: Scalable, Robust World Models for Continuous Control. In *International Conference on Learning Representations (ICLR)*, 2024.
- [19] Fangqi Zhu, Hongtao Wu, Song Guo, Yuxiao Liu, Chilam Cheang, and Tao Kong. IRASim: Learning Interactive Real-Robot Action Simulators. *arXiv preprint arXiv:2406.14540*, 2024.
- [20] Siyuan Zhou, Yilun Du, Jiaiben Chen, Yandong Li, Dit-Yan Yeung, and Chuang Gan. RoboDreamer: Learning Compositional World Models for Robot Imagination. *International Conference on Machine Learning (ICML)*, 2024.
- [21] Chongkai Gao, Haozhuo Zhang, Zhixuan Xu, Zhehao Cai, and Lin Shao. FLIP: Flow-Centric Generative Planning for General-Purpose Manipulation Tasks. *International Conference on Learning Representations (ICLR)*, 2025.
- [22] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-Sora: Democratizing Efficient Video Production for All. *arXiv preprint arXiv:2412.20404*, 2024.
- [23] Zheng Zhu, Xiaofeng Wang, Wangbo Zhao, Chen Min, Nianchen Deng, Min Dou, Yuqi Wang, Botian Shi, Kai Wang, Chi Zhang, et al. Is Sora a World Simulator? A Comprehensive Survey on General World Models and Beyond. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2024.
- [24] Younggyo Seo, Kimin Lee, Stephen L James, and Pieter Abbeel. Reinforcement Learning with Action-Free Pre-Training from Videos. In *International Conference on Machine Learning*, pages 19561–19579. PMLR, 2022.
- [25] Jialong Wu, Haoyu Ma, Chaoyi Deng, and Mingsheng Long. Pre-Training Contextualized World Models with In-the-Wild Videos for Reinforcement Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 36, 2024.
- [26] Jialong Wu, Shaofeng Yin, Ningya Feng, Xu He, Dong Li, Jianye Hao, and Mingsheng Long. iVideoGPT: Interactive VideoGPTs are Scalable World Models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2025.
- [27] Zhecheng Yuan, Zhengrong Xue, Bo Yuan, Xueqian Wang, Yi Wu, Yang Gao, and Huazhe Xu. Pre-Trained Image Encoder for Generalizable Visual Reinforcement Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [28] Sai Rajeswar, Pietro Mazzaglia, Tim Verbelen, Alexandre Piché, Bart Dhoedt, Aaron Courville, and Alexandre Lacoste. Mastering the unsupervised reinforcement learning benchmark from pixels. In *International Conference on Machine Learning*, pages 28598–28617. PMLR, 2023.
- [29] 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 (NeurIPS)*, 2019.
- [30] Yifan Wu, George Tucker, and Ofir Nachum. Behavior Regularized Offline Reinforcement Learning. *arXiv preprint arXiv:1911.11361*, 2019.
- [31] Ilya Kostrikov, Rob Fergus, Jonathan Tompson, and Ofir Nachum. Offline Reinforcement Learning with Fisher Divergence Critic Regularization. In *International Conference on Machine Learning (ICML)*, 2021.
- [32] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline Reinforcement Learning with Implicit Q-Learning. *International Conference on Learning Representations (ICLR)*, 2022.
- [33] Ikechukwu Uchendu, Ted Xiao, Yao Lu, Banghua Zhu, Mengyuan Yan, Joséphine Simon, Matthew Bennice, Chuyuan Fu, Cong Ma, Jiantao Jiao, et al. Jump-Start Reinforcement Learning. In *International Conference on Machine Learning (ICML)*, 2023.
- [34] Denis Yarats, David Brandfonbrener, Hao Liu, Michael Laskin, Pieter Abbeel, Alessandro Lazaric, and Lerrel Pinto. Don’t Change the Algorithm, Change the Data: Exploratory Data for Offline Reinforcement Learning. *arXiv preprint arXiv:2201.13425*, 2022.
- [35] Philip J Ball, Laura Smith, Ilya Kostrikov, and Sergey Levine. Efficient Online Reinforcement Learning with Offline Data. *International Conference on Machine Learning (ICML)*, 2023.
- [36] Qiyang Li, Jason Zhang, Dibya Ghosh, Amy Zhang, and Sergey Levine. Accelerating Exploration with Unlabeled Prior Data. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.

- [37] Ryan Julian, Benjamin Swanson, Gaurav S Sukhatme, Sergey Levine, Chelsea Finn, and Karol Hausman. Never Stop Learning: The Effectiveness of Fine-Tuning in Robotic Reinforcement Learning. In *Conference on Robot Learning (CoRL)*, 2020.
- [38] Dmitry Kalashnikov, Jacob Varley, Yevgen Chebotar, Benjamin Swanson, Rico Jonschkowski, Chelsea Finn, Sergey Levine, and Karol Hausman. MT-OPT: Continuous Multi-Task Robotic Reinforcement Learning at Scale. *arXiv preprint arXiv:2104.08212*, 2021.
- [39] Tianhe Yu, Aviral Kumar, Yevgen Chebotar, Karol Hausman, Sergey Levine, and Chelsea Finn. Conservative Data Sharing for Multi-Task Offline Reinforcement Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [40] Ignat Georgiev, Varun Giridhar, Nicklas Hansen, and Animesh Garg. PWM: Policy Learning with Large World Models. *arXiv preprint arXiv:2407.02466*, 2024.
- [41] Serkan Cabi, Sergio Gómez Colmenarejo, Alexander Novikov, Ksenia Konyushkova, Scott Reed, Rae Jeong, Konrad Zolna, Yusuf Aytar, David Budden, Mel Vecerik, et al. Scaling Data-Driven Robotics with Reward Sketching and Batch Reinforcement Learning. In *Robotics: Science and Systems (RSS)*, 2020.
- [42] Avi Singh, Larry Yang, Kristian Hartikainen, Chelsea Finn, and Sergey Levine. End-to-End Robotic Reinforcement Learning without Reward Engineering. *Robotics: Science and Systems (RSS)*, 2019.
- [43] Andrew Y Ng, Stuart Russell, et al. Algorithms for Inverse Reinforcement Learning. In *International Conference on Machine Learning (ICML)*, 2000.
- [44] Pieter Abbeel and Andrew Y Ng. Apprenticeship Learning via Inverse Reinforcement Learning. In *International Conference on Machine Learning (ICML)*, 2004.
- [45] Jonathan Ho and Stefano Ermon. Generative Adversarial Imitation Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016.
- [46] Tianhe Yu, Aviral Kumar, Yevgen Chebotar, Karol Hausman, Chelsea Finn, and Sergey Levine. How to Leverage Unlabeled Data in Offline Reinforcement Learning. In *International Conference on Machine Learning (ICML)*, 2022.
- [47] Adam Stooke, Kimin Lee, Pieter Abbeel, and Michael Laskin. Decoupling Representation Learning from Reinforcement Learning. In *International Conference on Machine Learning (ICML)*, 2021.
- [48] Rutav Shah and Vikash Kumar. RRL: Resnet as Representation for Reinforcement Learning. In *International Conference on Machine Learning (ICML)*, 2021.
- [49] Che Wang, Xufang Luo, Keith Ross, and Dongsheng Li. VRL3: A Data-Driven Framework for Visual Deep Reinforcement Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [50] Yanchao Sun, Shuang Ma, Ratnesh Madaan, Rogerio Bonatti, Furong Huang, and Ashish Kapoor. SMART: Self-supervised Multi-task preTrAining with contRol Transformers. In *International Conference on Learning Representations*, 2023.
- [51] Yanjie Ze, Nicklas Hansen, Yinbo Chen, Mohit Jain, and Xiaolong Wang. Visual Reinforcement Learning with Self-Supervised 3D Representations. *Robotics and Automation Letters*, 2023.
- [52] Dibya Ghosh, Chethan Anand Bhateja, and Sergey Levine. Reinforcement Learning from Passive Data via Latent Intentions. In *International Conference on Machine Learning (ICML)*, 2023.
- [53] Jialong Wu, Haoyu Ma, Chaoyi Deng, and Mingsheng Long. Pre-Training Contextualized World Models with In-the-Wild Videos for Reinforcement Learning. *Advances in Neural Information Processing Systems*, 36:39719–39743, 2023.
- [54] Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to Control: Learning Behaviors by Latent Imagination. *International Conference on Learning Representations (ICLR)*, 2020.
- [55] Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering Atari with Discrete World Models. *International Conference on Learning Representations (ICLR)*, 2021.
- [56] Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering Diverse Domains through World Models. *arXiv preprint arXiv:2301.04104*, 2023.
- [57] Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James Davidson. Learning Latent Dynamics for Planning from Pixels. In *International Conference on Machine Learning (ICML)*, 2019.
- [58] Benjamin Eysenbach, Vivek Myers, Sergey Levine, and Ruslan Salakhutdinov. Contrastive Representations Make Planning Easy. In *Advances in Neural Information Processing Systems Workshop (NeurIPS Workshop)*, 2023.

- [59] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is All You Need. *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [60] Albert Gu, Karan Goel, and Christopher Ré. Efficiently Modeling Long Sequences with Structured State Spaces. *International Conference on Learning Representations (ICLR)*, 2022.
- [61] Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvassy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The Faiss Library. *arXiv preprint arXiv:2401.08281*, 2024.
- [62] Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-Driven Exploration by Self-Supervised Prediction. In *International Conference on Machine Learning (ICML)*, 2017.
- [63] Nicklas Hansen, Zhecheng Yuan, Yanjie Ze, Tongzhou Mu, Aravind Rajeswaran, Hao Su, Huazhe Xu, and Xiaolong Wang. On Pre-Training for Visuo-Motor Control: Revisiting a Learning-from-Scratch Baseline. In *International Conference on Machine Learning (ICML)*, 2023.
- [64] German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual Lifelong Learning with Neural Networks: A Review. *Neural Networks*, 2019.
- [65] Khimya Khetarpal, Matthew Riemer, Irina Rish, and Doina Precup. Towards Continual Reinforcement Learning: A Review and Perspectives. *Journal of Artificial Intelligence Research*, 2022.
- [66] Arun Mallya and Svetlana Lazebnik. PackNet: Adding Multiple Tasks to a Single Network by Iterative Pruning. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [67] Bingshan Hu, Tianyue H Zhang, Nidhi Hegde, and Mark Schmidt. Optimistic Thompson Sampling-Based Algorithms for Episodic Reinforcement Learning. In *Uncertainty in Artificial Intelligence (UAI)*, 2023.
- [68] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust Region Policy Optimization. In *International Conference on Machine Learning (ICML)*, 2015.
- [69] Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by Random Network Distillation. *International Conference on Learning Representations*, 2019.
- [70] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is All You Need: Learning Skills without a Reward Function. In *International Conference on Learning Representations (ICLR)*, 2019.
- [71] Deepak Pathak, Dhiraj Gandhi, and Abhinav Gupta. Self-Supervised Exploration via Disagreement. In *International Conference on Machine Learning (ICML)*, 2019.
- [72] Hao Liu and Pieter Abbeel. Behavior From the Void: Unsupervised Active Pre-Training. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [73] Ramanan Sekar, Oleh Rybkin, Kostas Daniilidis, Pieter Abbeel, Danijar Hafner, and Deepak Pathak. Planning to Explore via Self-Supervised World Models. In *International Conference on Machine Learning (ICML)*, 2020.
- [74] Hao Liu and Pieter Abbeel. APS: Active Pretraining with Successor Features. In *International Conference on Machine Learning (ICML)*, 2021.
- [75] Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Reinforcement Learning with Prototypical Representations. In *International Conference on Machine Learning (ICML)*, 2021.
- [76] Michael Laskin, Denis Yarats, Hao Liu, Kimin Lee, Albert Zhan, Kevin Lu, Catherine Cang, Lerrel Pinto, and Pieter Abbeel. URLB: Unsupervised Reinforcement Learning Benchmark. *Neural Information Processing Systems Track on Datasets and Benchmarks*, 2021.
- [77] Pietro Mazzaglia, Tim Verbelen, Bart Dhoedt, Alexandre Lacoste, and Sai Rajeswar. Choreographer: Learning and adapting skills in imagination. In *International Conference on Learning Representations (ICLR)*, 2023.
- [78] Yingchen Xu, Jack Parker-Holder, Aldo Pacchiano, Philip Ball, Oleh Rybkin, S Roberts, Tim Rocktäschel, and Edward Grefenstette. Learning General World Models in a Handful of Reward-Free Deployments. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [79] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster Level in StarCraft II using Multi-Agent Reinforcement Learning. *Nature*, 2019.
- [80] OpenAI: Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob McGrew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, et al. Learning Dexterous In-Hand Manipulation. *The International Journal of Robotics Research (IJRR)*, 2020.

- [81] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al. A Generalist Agent. *Transactions on Machine Learning Research (TMLR)*, 2022.
- [82] Anthony Brohan, Noah Brown, Justice Carbalaj, 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. *Robotics: Science and Systems (RSS)*, 2023.
- [83] 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. *Robotics: Science and Systems (RSS)*, 2024.
- [84] Yi Zhao, Le Chen, Jan Schneider, Quankai Gao, Juho Kannala, Bernhard Schölkopf, Joni Pajarinen, and Dieter Büchler. RP1M: A Large-Scale Motion Dataset for Piano Playing with Bi-Manual Dexterous Robot Hands. *Conference on Robot Learning (CoRL)*, 2024.
- [85] Anthony Brohan, Noah Brown, Justice Carbalaj, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control. In *Conference on Robot Learning (CoRL)*, 2023.
- [86] Abby O'Neill, Abdul Rehman, Abhinav Gupta, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, et al. Open X-Embodiment: Robotic Learning Datasets and RT-X Models. In *International Conference on Robotics and Automation (ICRA)*, 2023.
- [87] 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. In *Robotics: Science and Systems (RSS)*, 2024.
- [88] David Ha and Jürgen Schmidhuber. Recurrent World Models Facilitate Policy Evolution. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [89] Vincent Micheli, Eloi Alonso, and François Fleuret. Transformers are Sample-Efficient World Models. In *International Conference on Learning Representations (ICLR)*, 2023.
- [90] Eloi Alonso, Adam Jelleby, Vincent Micheli, Anssi Kanervisto, Amos Storkey, Tim Pearce, and François Fleuret. Diffusion for World Modeling: Visual Details Matter in Atari. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- [91] Aidan Scannell, Mohammadreza Nakhaei, Kalle Kujanpää, Yi Zhao, Kevin Luck, Arno Solin, and Joni Pajarinen. Discrete Codebook World Models for Continuous Control. In *International Conference on Learning Representations (ICLR)*, 2025.
- [92] Anthony Hu, Lloyd Russell, Hudson Yeo, Zak Murez, George Fedoseev, Alex Kendall, Jamie Shotton, and Gianluca Corrado. GAIA-1: A Generative World Model for Autonomous Driving. *arXiv preprint arXiv:2309.17080*, 2023.
- [93] Tim Pearce, Tabish Rashid, Dave Bignell, Raluca Georgescu, Sam Devlin, and Katja Hofmann. Scaling laws for pre-training agents and world models. *arXiv preprint arXiv:2411.04434*, 2024.
- [94] Niket Agarwal, Arslan Ali, Maciej Bala, Yogesh Balaji, Erik Barker, Tiffany Cai, Prithvijit Chattopadhyay, Yongxin Chen, Yin Cui, Yifan Ding, et al. Cosmos World Foundation Model Platform for Physical AI. In *arXiv preprint arXiv:2501.03575*, 2025.
- [95] Yann LeCun, Yoshua Bengio, et al. Convolutional Networks for Images, Speech, and Time Series. *The Handbook of Brain Theory and Neural Networks*, 1995.
- [96] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer Normalization. *arXiv preprint arXiv:1607.06450*, 2016.

Appendices

A More Results	16
B Theoretical Analysis	17
C More Related Work	19
D Limitations	19
E Impact Statement	19
F Compute Resources	20
G Implementation Details	20
H Algorithm	22
I Full Results	23
J Hyperparameters	26
K Task Visualization	27

A More Results

Comparison with Imitation Learning Baseline To demonstrate the mix-quality property of the non-curated dataset, we compare GSA with Diffusion Policy, a widely used imitation learning approach by modeling the agent with diffusion models. From App. A, we can see that due to the dataset consisting of non-expert data, the diffusion policy fails to demonstrate satisfactory results, while GSA can effectively utilize the offline data.

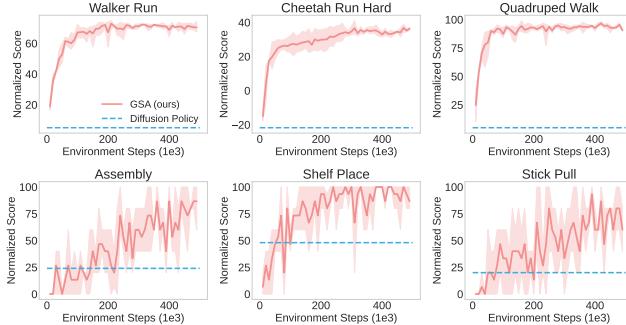


Figure 9: Comparison with Diffusion Policy. GSA can effectively handle non-curated offline data while the imitation learning baseline fails.

Full Results of Comparison with iVideoGPT We compare with other model-based approaches on tasks used in iVideoGPT [26]. We show that GSA outperforms the baselines without using reward shaping and pre-filling the replay buffer with demonstrations. This highlights that although non-curated, the offline data can clearly boost RL training, and GSA can effectively use the information in the data.

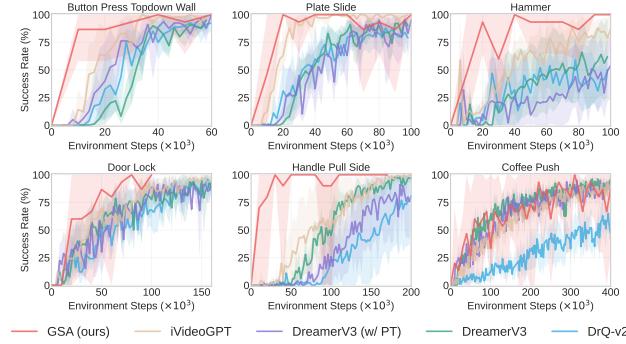


Figure 10: Comparsion with model-based approaches for leveraging offline data.

Role of Each Component We offer inter-quartile mean(IQM) and optimality gap for the ablation study of the role of each proposed component. Together with the retrieval-based experience rehearsal and execution guidance, a pre-trained generalist world model boosts RL performance on a wide range of tasks.

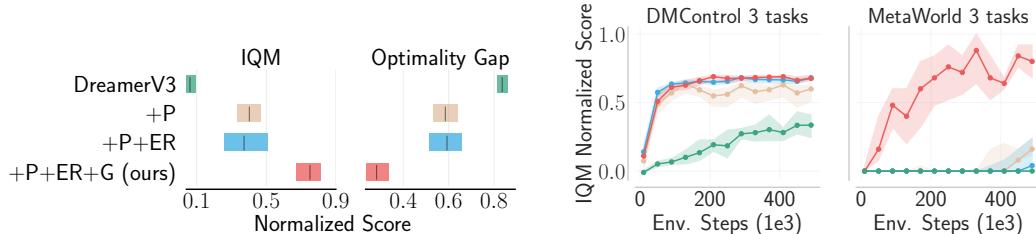


Figure 11: Ablation study on the role of each component. “P” represents world model pretraining, “ER” means experience rehearsal, and “G” represents execution guidance. Together with the proposed retrieval-based experience rehearsal and execution guidance, world model pre-training boosts RL performance on a wide range of tasks.

B Theoretical Analysis

In this section, we give a theoretical analysis of the main conclusions in our paper.

B.1 Proof of the Benefits of Experience Retrieval

Proposition 1. *Experience retrieval reduces distribution shift during online fine-tuning, compared to using the full offline dataset directly, in the sense that*

$$\mathbb{E}_{s \sim p_{\text{retrieved}}, s_{\text{on}} \sim p_{\text{on}}} [\|s - s_{\text{on}}\|_2] < \mathbb{E}_{s \sim p_{\text{off}}, s_{\text{on}} \sim p_{\text{on}}} [\|s - s_{\text{on}}\|_2]. \quad (5)$$

Proof. Let $p_{\text{off}}(s)$, $p_{\text{on}}(s)$, and $p_{\text{retrieved}}(s)$ denote the state distributions of the non-curated offline dataset \mathcal{D}_{off} , the online dataset \mathcal{D}_{on} , and the retrieved dataset $\mathcal{D}_{\text{retrieved}} \subset \mathcal{D}_{\text{off}}$, respectively. We simplify the notation as

$$\mathbb{E}_{s \sim p, s_{\text{on}} \sim p_{\text{on}}} [\|s - s_{\text{on}}\|_2] \quad \text{as} \quad \mathbb{E}_{s \sim p} [d(s, s_{\text{on}})].$$

Since $\mathcal{D}_{\text{retrieved}} \subset \mathcal{D}_{\text{off}}$, the distribution $p_{\text{off}}(s)$ can be expressed as a mixture distribution:

$$p_{\text{off}}(s) = \alpha \cdot p_{\text{retrieved}}(s) + (1 - \alpha) \cdot p_{\text{rest}}(s),$$

where $p_{\text{rest}}(s)$ is the distribution over the remaining offline data, and $\alpha = \frac{|\mathcal{D}_{\text{retrieved}}|}{|\mathcal{D}_{\text{off}}|}$ denotes the fraction of samples in the retrieved dataset.

The expected total variation for the mixture distribution decomposes as:

$$\mathbb{E}_{s \sim p_{\text{off}}} [d(s, s_{\text{on}})] = \alpha \cdot \mathbb{E}_{s \sim p_{\text{retrieved}}} [d(s, s_{\text{on}})] + (1 - \alpha) \cdot \mathbb{E}_{s \sim p_{\text{rest}}} [d(s, s_{\text{on}})]. \quad (6)$$

Assume that $\mathcal{D}_{\text{retrieved}}$ is constructed by selecting states such that $\|s_{\text{retrieved}} - s_{\text{on}}\| < \epsilon$, for some small $\epsilon > 0$. Consequently, states in $\mathcal{D}_{\text{rest}}$ satisfy $\|s_{\text{rest}} - s_{\text{on}}\| \geq \epsilon$. This construction implies the following bounds:

$$\mathbb{E}_{s \sim p_{\text{retrieved}}} [d(s, s_{\text{on}})] < \epsilon', \quad (7)$$

$$\mathbb{E}_{s \sim p_{\text{rest}}} [d(s, s_{\text{on}})] \geq \epsilon', \quad (8)$$

for some $\epsilon' > 0$. Therefore, it follows that

$$\mathbb{E}_{s \sim p_{\text{rest}}} [d(s, s_{\text{on}})] > \mathbb{E}_{s \sim p_{\text{retrieved}}} [d(s, s_{\text{on}})].$$

Substituting into Equation (6) yields:

$$\begin{aligned} \mathbb{E}_{s \sim p_{\text{off}}} [d(s, s_{\text{on}})] &= \alpha \cdot \mathbb{E}_{s \sim p_{\text{retrieved}}} [d(s, s_{\text{on}})] + (1 - \alpha) \cdot \mathbb{E}_{s \sim p_{\text{rest}}} [d(s, s_{\text{on}})] \\ &> \alpha \cdot \mathbb{E}_{s \sim p_{\text{retrieved}}} [d(s, s_{\text{on}})] + (1 - \alpha) \cdot \mathbb{E}_{s \sim p_{\text{retrieved}}} [d(s, s_{\text{on}})] \\ &= \mathbb{E}_{s \sim p_{\text{retrieved}}} [d(s, s_{\text{on}})]. \end{aligned}$$

Thus, the expected total variation between the retrieved data and online data is strictly smaller than that between the full offline data and online data. \square

Proposition 2. *Experience retrieval helps prevent catastrophic forgetting during online fine-tuning.*

Definition 1 (Catastrophic Forgetting due to Data Distribution Shift). *Catastrophic forgetting occurs when a neural network, after training on a new data distribution, experiences a significant performance drop on previously learned tasks due to the overwriting of representations from earlier distributions, caused by biased parameter updates towards the new distribution.*

Proof. Following the previous notations, let \mathcal{D}_{on} and $\mathcal{D}_{\text{retrieved}}$ denote the online dataset and the retrieved offline dataset, respectively. The objective in Eq. (1) can be written as:

$$\begin{aligned} \mathcal{L}_{\text{mixed}}(\theta) &= \mathcal{L}_{\text{on}}(\theta) + \lambda \cdot \mathcal{L}_{\text{retrieved}}(\theta) \\ &= \mathbb{E}_{p_{\theta}, q_{\theta}, (o, a) \sim \mathcal{D}_{\text{on}}} \left[\sum_{t=1}^T -\ln p_{\theta}(o_t | z_t, h_t) + \beta \cdot \text{KL}(q_{\theta}(z_t | h_t, o_t) \| p_{\theta}(z_t | h_t)) \right] \\ &\quad + \lambda \cdot \mathbb{E}_{p_{\theta}, q_{\theta}, (o, a) \sim \mathcal{D}_{\text{retrieved}}} \left[\sum_{t=1}^T -\ln p_{\theta}(o_t | z_t, h_t) + \beta \cdot \text{KL}(q_{\theta}(z_t | h_t, o_t) \| p_{\theta}(z_t | h_t)) \right]. \end{aligned}$$

Assuming the λ is a monotonic function of $\alpha = \frac{|\mathcal{D}_{\text{retrieved}}|}{|\mathcal{D}_{\text{off}}|}$ and $\lambda > 0$, since $\mathcal{D}_{\text{retrieved}} \subset \mathcal{D}_{\text{off}}$, the term $\mathcal{L}_{\text{retrieved}}(\theta)$ acts as a regularizer during online updates, constraining parameter changes on \mathcal{D}_{on} in a way that preserves performance on the retrieved offline distribution $p_{\text{retrieved}}$. This mitigates the risk of catastrophic forgetting by anchoring the model to previously seen data. \square

B.2 Proof of Improved Performance with Execution Guidance

Proposition 3 (Performance Improvement via Execution Guidance). *Let π^e denote an exploration policy and π^g a guide policy obtained via imitation learning. Assume that π^g outperforms π^e at the early stage of training, that is, there exists a constant $\delta > 0$ such that for all states s in the support of the state visitation distribution:*

$$\mathbb{E}_{a \sim \pi^g(\cdot|s)}[A_{\pi^e}(s, a)] \geq \delta \quad (9)$$

Let $\tilde{\pi}$ be an α -coupled policy (execution guidance) derived from π^e and π^g , as defined in Definition 1 of [68], such that:

$$P(a \neq a_g | s) \leq \alpha, \quad (10)$$

where $a \sim \pi^e(\cdot|s)$ and $a_g \sim \pi^g(\cdot|s)$.

Then, the performance of the α -coupled policy $\tilde{\pi}$ exceeds that of the exploration policy π^e by at least:

$$\eta(\tilde{\pi}) - \eta(\pi^e) \geq \frac{\alpha \cdot \delta}{1 - \gamma}, \quad (11)$$

where $\gamma \in [0, 1]$ is the discount factor.

Proof. By Lemma 1 in Trust Region Policy Optimization (TRPO) [68], the difference in policy performance can be expressed as:

$$\eta(\tilde{\pi}) = \eta(\pi^e) + \mathbb{E}_{\tau \sim \tilde{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t A_{\pi^e}(s_t, a_t) \right] \quad (12)$$

Define the expected advantage at state s as:

$$\bar{A}_{\pi^e}(s) = \mathbb{E}_{a \sim \tilde{\pi}(\cdot|s)}[A_{\pi^e}(s, a)]$$

Given the definition of an α -coupled policy, we can decompose this expectation:

$$\bar{A}_{\pi^e}(s) = (1 - \alpha) \cdot \mathbb{E}_{a \sim \pi^e(\cdot|s)}[A_{\pi^e}(s, a)] + \alpha \cdot \mathbb{E}_{a \sim \pi^g(\cdot|s)}[A_{\pi^e}(s, a)] \quad (13)$$

Since $\mathbb{E}_{a \sim \pi^e(\cdot|s)}[A_{\pi^e}(s, a)] = 0$ (by definition of the advantage function), we have:

$$\begin{aligned} \bar{A}_{\pi^e}(s) &= \alpha \cdot \mathbb{E}_{a \sim \pi^g(\cdot|s)}[A_{\pi^e}(s, a)] \\ &\geq \alpha \cdot \delta \end{aligned} \quad (14)$$

The policy performance difference can then be bounded as:

$$\begin{aligned} \eta(\tilde{\pi}) - \eta(\pi^e) &= \mathbb{E}_{\tau \sim \tilde{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t \bar{A}_{\pi^e}(s_t) \right] \\ &= \sum_s \rho_{\tilde{\pi}}(s) \bar{A}_{\pi^e}(s) \\ &\geq \sum_s \rho_{\tilde{\pi}}(s) \cdot \alpha \cdot \delta \\ &= \alpha \cdot \delta \cdot \sum_s \rho_{\tilde{\pi}}(s), \end{aligned} \quad (15)$$

where $\rho_{\tilde{\pi}}(s)$ is the discounted state visitation frequency under $\tilde{\pi}$.

Since $\sum_s \rho_{\tilde{\pi}}(s) = \frac{1}{1-\gamma}$ (the sum of discounted state visitation frequencies equals the effective horizon length), we have:

$$\eta(\tilde{\pi}) - \eta(\pi^e) \geq \frac{\alpha \cdot \delta}{1 - \gamma} \quad (16)$$

This establishes that the performance improvement of the α -coupled policy $\tilde{\pi}$ over the exploration policy π^e is at least $\frac{\alpha \cdot \delta}{1 - \gamma}$, which is strictly positive when $\alpha > 0$, $\delta > 0$, and $\gamma < 1$. \square

C More Related Work

Unsupervised RL In unsupervised RL, an agent explores the environment based on intrinsic motivations, and the models’ parameters are initialized during this self-motivated exploration stage, aiming for fast downstream task learning [28, 62, 69–78]. Our problem setting differs from unsupervised RL in several ways: i) Unsupervised RL interacts with the environment actively while we leverage *static* offline datasets, ii) unsupervised RL gives a specific focus on designing different intrinsic rewards, while our setting focuses on improving sample efficiency by leveraging unlabeled datasets.

Generalist Agents RL methods usually perform well on a single task [79, 80], however, this contrasts with humans who can perform multiple tasks well. Recent works have proposed generalist agents that master a diverse set of tasks with a single agent [81–84]. These methods typically resort to scalable models and large datasets and are trained via imitation learning [82, 85–87]. In contrast, we train a generalist world model and use it to boost RL performance for multiple tasks and embodiments.

World models World models learn to predict future observations or states based on historical information. World models have been widely investigated in online model-based RL [54, 88–91]. Recently, the community has started investigating scaling world models [88], for example, Hu et al. [92], Pearce et al. [93], Wu et al. [26], Agarwal et al. [94] train world models with Diffusion Models or Transformers. However, these models are usually trained on demonstration data. In contrast, we explore the offline-to-online RL setting – closely fitting the pre-train and then fine-tune paradigm – and we focus on leveraging reward-free and multi-embodiment data to increase the amount of available data for pre-training. We further identify the distributional shift issue when fine-tuning the pre-trained world model and mitigate the issue by proposing experience rehearsal and execution guidance.

D Limitations

Although demonstrating strong performance on a diverse set of tasks, our method has the following limitations. 1) The world model architecture used in our paper is the recurrent state space model. This model is built upon RNN, which can be limited for scaling. This can be mitigated by using a Transformer and a diffusion-based world model. However, we note that the main conclusion of this paper should still be valid. 2) We do not thoroughly discuss the generalization ability of the pre-trained world model. With DMControl tasks, our method shows a promising trend in generalizing to unseen tasks. However, generalization to new embodiments or novel configurations is still challenging, which requires even diverse training data. 3) The non-curated offline data used in our paper, although lifting several key assumptions in previous offline-to-online RL, is still in-domain data, i.e., our current method is not able to leverage the vast in-the-wild data. A promising direction is to combine in-the-wild data for pre-training as in [26] and the domain-specific “in-house” data (as used in our paper) for post-training. 4) We only conduct experiments in the simulator. Considering the sample efficiency of our proposed method, it could be promising to conduct experiments on real-world applications.

E Impact Statement

This paper contributes to the field of reinforcement learning (RL), with potential applications including robotics and autonomous machines. While our methods hold promise for advancing technology, they could also be applied in ways that raise ethical concerns, such as in autonomous machines exploring the world and making decisions on their own. However, the specific societal impacts of our work are

broad and varied, and we believe a detailed discussion of potential negative uses is beyond the scope of this paper. We encourage a broader dialogue on the ethical use of RL technology and its regulation to prevent misuse.

F Compute Resources

We conduct all experiments on clusters equipped with AMD MI250X GPUs, 64-core AMD EPYC "Trento" CPUs, and 64 GBs DDR4 memory. For pre-training, it takes \sim 48 GPU hours for 150k steps. For fine-tuning, it tasks \sim 8 GPU hours per run for 150K environment steps. Note that due to AMD GPUs not supporting hardware rendering, the training time should be longer than using Nvidia GPUs. To reproduce the GSA's results in Fig. 3, it roughly takes $8\text{ h} * 72\text{ tasks} * 3\text{ seeds} = 1728$ GPU hours.

G Implementation Details

G.1 Behavior Cloning

The Behavior Cloning methods used in both the execution guidance of GSA and JSRL-BC are the same. We use a four-layer convolutional neural network [95] with kernel depth [32, 64, 128, 256] following a three-layer MLPs with LayerNorm [96] after all linear layers.

We list the adopted encoder and actor architectures for reference.

```

1 class Encoder(nn.Module):
2     def __init__(self, obs_shape):
3         super().__init__()
4         assert obs_shape == (9, 64, 64), f'obs_shape is {(obs_shape)}, but'
5         expect (9, 64, 64)' # inputs shape
6
7         self.repr_dim = (32 * 8) * 2 * 2
8         _input_channel = 9
9
10        self.convnet = nn.Sequential(
11            nn.Conv2d(_input_channel, 32, 4, stride=2), # [B, 32, 31, 31]
12            nn.ELU(),
13            nn.Conv2d(32, 32*2, 4, stride=2), #[B, 64, 14, 14]
14            nn.ELU(),
15            nn.Conv2d(32*2, 32*4, 4, stride=2), #[B, 128, 6, 6]
16            nn.ELU(),
17            nn.Conv2d(32*4, 32*8, 4, stride=2), #[B, 256, 2, 2]
18            nn.ELU())
19        self.apply(utils.weight_init)
20
21    def forward(self, obs):
22        B, C, H, W = obs.shape
23
24        obs = obs / 255.0 - 0.5
25        h = self.convnet(obs)
26        # reshape to [B, -1]
27        h = h.view(B, -1)
28        return h

```

```

11             nn.Linear(hidden_dim, action_shape[0]))
12         )
13     self.apply(utils.weight_init)
14
15     def forward(self, obs, std):
16         h = self.trunk(obs)
17         return self.policy(h)

```

G.2 JSRL+BC

Jump-start RL [33] is proposed as an offline-to-online RL method. It includes two policies, a prior policy $\pi_{\theta_1}(a|s)$ and a behavior policy $\pi_{\theta_2}(a|s)$, where the prior policy is trained via offline RL methods and the behavior policy is updated during the online learning stage. However, offline RL requires the offline dataset to include rewards for the target task. To extract behavior policy from the offline dataset, we use the BC agent described above as the prior policy. During online training, in each episode, we randomly sample the rollout horizon h of the prior policy from a pre-defined array `np.arange(0, 101, 10)`. We then execute the prior policy for h steps and switch to the behavior policy until the end of an episode.

G.3 ExPLORe

For the ExPLORe baseline, we follow the original training code ¹. We sweep over several design choices: i) kernel size of the linear layer used in the RND and reward models: [256 (default), 512]; ii) initial temperature value: [0.1 (default), 1.0]; iii) whether to use LayerNorm Layer (no by default); iv) learning rate: [1e-4, 3e-4 (default)]. However, we fail to obtain satisfactory performance. There are several potential reasons: i) the parameters used in the ExPLORe paper are tuned specifically to their setting, where manipulation tasks and near-expert trajectories are used; ii) the coefficient term of the RND value needs to be tuned carefully for different tasks and the reward should also be properly normalized.

To achieve reasonable performance and eliminate the performance gap caused by implementation-level details, we make the following modifications: i) we replace the RND module with ensembles to calculate uncertainty; ii) the reward function share the latent space with actor and critic.

¹Source code of ExPLORe <https://github.com/facebookresearch/ExPLORe>

H Algorithm

The full algorithm is described in Alg. 1.

Algorithm 1 Efficient RL by Guiding Generalist World Models with Non-Curated Offline Data

Require: Non-curated offline data \mathcal{D}_{off} , Online data $\mathcal{D}_{\text{on}} \leftarrow \emptyset$, Retrieval data $\mathcal{D}_{\text{retrieval}} \leftarrow \emptyset$
 World model $f_\theta, q_\theta, p_\theta, d_\theta$
 Policy $\pi_{\phi_{\text{RL}}}, \pi_{\phi_{\text{BC}}}$, Value function v_ϕ and Reward r_ξ .

// Task-Agnostic World Model Pre-Training

for num. pre-train steps **do**

- Randomly sample mini-batch $\mathcal{B}_{\text{off}} : \{o_t, a_t, o_{t+1}\}_{t=0}^T$ from \mathcal{D}_{off} .
- Update world model $f_\theta, q_\theta, p_\theta, d_\theta$ by minimizing Eq. (1) on sampled batch \mathcal{B} .

end for

// Task-Specific Training

// Experience Retrieval

Collect one initial observation o_{on}^0 from the environment.
 Compute the visual similarity between o_{on} and initial observations of trajectories o_{off} in \mathcal{D}_{off} using Eq. (4).
 Select R trajectories according to Eq. (4) and fill $\mathcal{D}_{\text{retrieval}}$.

// Behavior Cloning Policy Training

for num. bc updates **do**

- Randomly sample mini-batch $\mathcal{B}_{\text{retrieval}} : \{o_i, a_i\}_{i=0}^N$ from $\mathcal{D}_{\text{retrieval}}$.
- Update $\pi_{\phi_{\text{BC}}}$ by minimizing $-\frac{1}{N} \sum_{i=0}^N \log \pi_{\phi_{\text{BC}}}(a_i | o_i)$.

end for

// Task-Specific RL Fine-Tuning

for num. episodes **do**

// Collect Data

Decide whether to use $\pi_{\phi_{\text{BC}}}$ according to the predefined schedule.

if Select $\pi_{\phi_{\text{BC}}}$ **then**

- Randomly select the starting time step k and the rollout horizon H .

end if

$t \leftarrow 0$

while $t \leq$ episode length **do**

- $a_t = \pi_{\phi_{\text{BC}}}(a_t | o_t)$ if Use $\pi_{\phi_{\text{BC}}}$ and $k \leq t \leq H$ else $a_t = \pi_{\phi_{\text{RL}}}(a_t | o_t)$.
- Interact the environment with a_t . Store $\{o_t, a_t, r_t, o_{t+1}\}$ to \mathcal{D}_{on} .
- $t \leftarrow t + 1$

end while

// Update Models

for num. grad steps **do**

- Randomly sample mini-batch $\mathcal{B}_{\text{on}} : \{o_t, a_t, r_t, o_{t+1}\}_{t=0}^T$ from \mathcal{D}_{on} and $\mathcal{B}_{\text{retrieval}} : \{o_t, a_t, r_t, o_{t+1}\}_{t=0}^T$ from $\mathcal{D}_{\text{retrieval}}$.
- Update world model $f_\theta, q_\theta, p_\theta, d_\theta$ by minimizing Eq. (1) on sampled batch $\{\mathcal{B}_{\text{on}}, \mathcal{B}_{\text{retrieval}}\}$.
- Update r_ξ by minimizing $-\frac{1}{N} \sum_{i=0}^N \log p_\xi(r_i | s_i)$ on \mathcal{B}_{on} . $\triangleleft s_t = [h_t, z_t]$

// Update policy and value function

Generate imaginary trajectories $\tilde{\tau} = \{s_t, a_t, s_{t+1}\}_{t=0}^T$ by rolling out h_θ, p_θ with $\pi_{\phi_{\text{RL}}}$.

Update policy $\pi_{\phi_{\text{RL}}}$ and value function v_ϕ with Eq. (3).

end for

end for

I Full Results

In [Table 2](#) and [Table 3](#), we list the success rate of 50 Meta-World benchmark tasks with pixel inputs. In [Table 4](#), we list the episodic return of DMControl of 22 tasks. We compare GSA at 150k samples with two widely used baselines Dreamer v3 and DrQ v2 at both 150k samples and 1M samples. We mark the best result with a bold font at 150k samples and use underlining to mark the highest score overall.

I.1 Meta-World Benchmark

Table 2: Success rate of Meta-World benchmark with pixel inputs.

Tasks	Dreamer v3 @ 1M	DrQ v2 @ 1M	Dreamer v3 @ 150k	DrQ v2 @ 150k	GSA @ 150k
Assembly	0.0	0.0	0.0	0.0	0.2
Basketball	0.0	<u>0.97</u>	0.0	0.0	0.4
Bin Picking	0.0	<u>0.93</u>	0.0	0.33	0.8
Box Close	0.13	<u>0.9</u>	0.0	0.0	0.9
Button Press	<u>1.0</u>	0.7	0.47	0.13	0.9
Button Press Topdown	<u>1.0</u>	<u>1.0</u>	0.33	0.17	1.0
Button Press Topdown Wall	<u>1.0</u>	<u>1.0</u>	0.73	0.63	1.0
Button Press Wall	<u>1.0</u>	<u>1.0</u>	0.93	0.77	1.0
Coffee Button	1.0	1.0	1.0	1.0	1.0
Coffee Pull	0.6	<u>0.8</u>	0.0	0.6	0.6
Coffee Push	0.67	<u>0.77</u>	0.13	0.2	0.7
Dial Turn	<u>0.67</u>	0.43	0.13	0.17	0.67
Disassemble	0.0	0.0	0.0	0.0	0.0
Door Close	-	-	-	-	1.0
Door Lock	<u>1.0</u>	0.93	0.6	0.97	0.9
Door Open	<u>1.0</u>	0.97	0.0	0.0	0.8
Door Unlock	1.0	1.0	1.0	0.63	0.8
Drawer Close	0.93	1.0	0.93	1.0	0.9
Drawer Open	0.67	0.33	0.13	0.33	1.0
Faucet Open	1.0	1.0	0.47	0.33	1.0
Faucet Close	0.87	1.0	1.0	1.0	0.8
Hammer	1.0	1.0	0.07	0.4	1.0
Hand Insert	0.07	<u>0.57</u>	0.0	0.1	0.4
Handle Press Side	1.0	1.0	1.0	1.0	1.0
Handle Press	1.0	1.0	0.93	0.97	1.0
Handle Pull Side	0.67	1.0	0.67	0.6	1.0
Handle Pull	0.67	0.6	0.33	0.6	1.0
Lever Pull	0.73	<u>0.83</u>	0.0	0.33	0.8

More results see [Table 3](#)

Table 3: Success rate of Meta-World benchmark with pixel inputs (Cont.).

Tasks	Dreamer v3 @ 1M	DrQ v2 @ 1M	Dreamer v3 @ 150k	DrQ v2 @ 150k	GSA (ours) @ 150k
Peg Insert Side	1.0	1.0	0.0	0.27	<u>1.0</u>
Peg Unplug Side	<u>0.93</u>	0.9	0.53	0.5	0.8
Pick Out of Hole	0.0	0.27	0.0	0.0	<u>0.3</u>
Pick Place Wall	0.2	0.17	0.0	0.0	<u>0.5</u>
Pick Place	<u>0.67</u>	<u>0.67</u>	0.0	0.0	0.2
Plate Slide Back Side	1.0	1.0	0.93	1.0	<u>1.0</u>
Plate Slide Back	1.0	1.0	0.8	0.97	<u>1.0</u>
Plate Slide Side	<u>1.0</u>	0.9	0.73	0.5	0.5
Plate Slide	1.0	1.0	0.93	1.0	<u>1.0</u>
Push Back	<u>0.33</u>	<u>0.33</u>	0.0	0.0	0.2
Push Wall	0.33	0.57	0.0	0.0	<u>0.9</u>
Push	0.26	<u>0.93</u>	0.0	0.13	0.7
Reach	<u>0.87</u>	0.73	0.67	0.43	0.3
Reach Wall	<u>1.0</u>	0.87	0.53	0.7	<u>0.9</u>
Shelf Place	0.4	0.43	0.0	0.0	<u>0.87</u>
Soccer	0.6	0.3	0.13	0.13	<u>0.67</u>
Stick Push	0.0	0.07	0.0	0.0	0.4
Stick Pull	0.0	0.33	0.0	0.0	<u>0.67</u>
Sweep Into	0.87	<u>1.0</u>	0.0	0.87	0.9
Sweep	0.0	<u>0.73</u>	0.0	0.3	0.6
Window Close	1.0	1.0	0.93	1.0	0.8
Window Open	1.0	0.97	0.6	1.0	0.9
Mean	<u>0.900</u>	0.753	0.360	0.442	0.750
Medium	0.870	<u>0.900</u>	0.130	0.330	0.800

I.2 DMControl Benchmark

Table 4: Episodic return of DMControl benchmark with pixel inputs.

Tasks	Dreamer v3 @ 500k	DrQ v2 @ 500k	Dreamer v3 @ 150k	DrQ v2 @ 150k	GSA(ours) @ 150k
CartPole Balance	994.3	992.3	955.8	983.3	995.0
Acrobot Swingup	<u>222.1</u>	30.3	85.2	20.8	50.3
Acrobot Swingup Sparse	2.5	1.17	1.7	1.5	26.9
Acrobot Swingup Hard	-0.2	0.3	2.0	0.4	10.7
Walker Stand	965.7	947.6	946.2	742.9	969.4
Walker Walk	949.2	797.8	808.9	280.1	959.1
Walker Run	616.6	299.3	224.4	143.0	728.0
Walker Backflip	293.6	96.7	128.2	91.7	306.0
Walker Walk Backward	<u>942.9</u>	744.3	625.9	470.9	863.6
Walker Walk Hard	-2.1	-9.5	-4.7	-17.1	878.3
Walker Run Backward	<u>363.8</u>	246.0	229.4	167.4	349.3
Cheetah Run	<u>843.7</u>	338.1	621.4	251.2	526.1
Cheetah Run Front	<u>473.8</u>	202.4	143.1	108.4	360.5
Cheetah Run Back	<u>657.4</u>	294.4	407.6	171.2	446.0
Cheetah Run Backwards	<u>693.8</u>	384.3	626.6	335.6	542.2
Cheetah Jump	597.0	535.6	200.8	251.8	634.1
Quadruped Walk	369.3	258.1	145.2	76.5	933.6
Quadruped Stand	746.0	442.2	227.2	318.9	936.4
Quadruped Run	328.1	296.5	183.0	102.8	802.5
Quadruped Jump	689.6	478.3	168.3	190.5	813.5
Quadruped Roll	663.9	446.0	207.9	126.2	970.8
Quadruped Roll Fast	508.8	366.9	124.8	164.7	782.0
Mean	541.81	372.23	320.86	226.49	631.10
Medium	606.8	318.70	204.35	166.05	755.0

J Hyperparameters

In this section, we list important hyperparameters used in GSA.

Table 5: Hyperparameters used in GSA.

Hyperparameter	Value
Pre-training	
Stacked images	1
Pretrain steps	200,000
Batch size	16
Sequence length	64
Replay buffer capacity	Unlimited
Replay sampling strategy	Uniform
RSSM	
Hidden dimension	12288
Deterministic dimension	1536
Stochastic dimension	32 * 96
Block number	8
Layer Norm	True
CNN channels	[96, 192, 384, 768]
Activation function	SiLU
Optimizer	
Optimizer	Adam
Learning rate	1e-4
Weight decay	1e-6
Eps	1e-5
Gradient clip	100
Fine-tuning	
Warm-up frames	15000
Execution Guidance Schedule	linear(1,0,50000) for DMControl linear(1,0,1,150000) for Meta-Wolrd
Action repeat	2
Offline data mix ratio	0.25
Discount	0.99
Discount lambda	0.95
MLPs	[512, 512, 512]
MLPs activation	SiLU
Actor critic learning rate	8e-5
Actor entropy coef	1e-4
Target critic update fraction	0.02
Imagine horizon	16

K Task Visualization

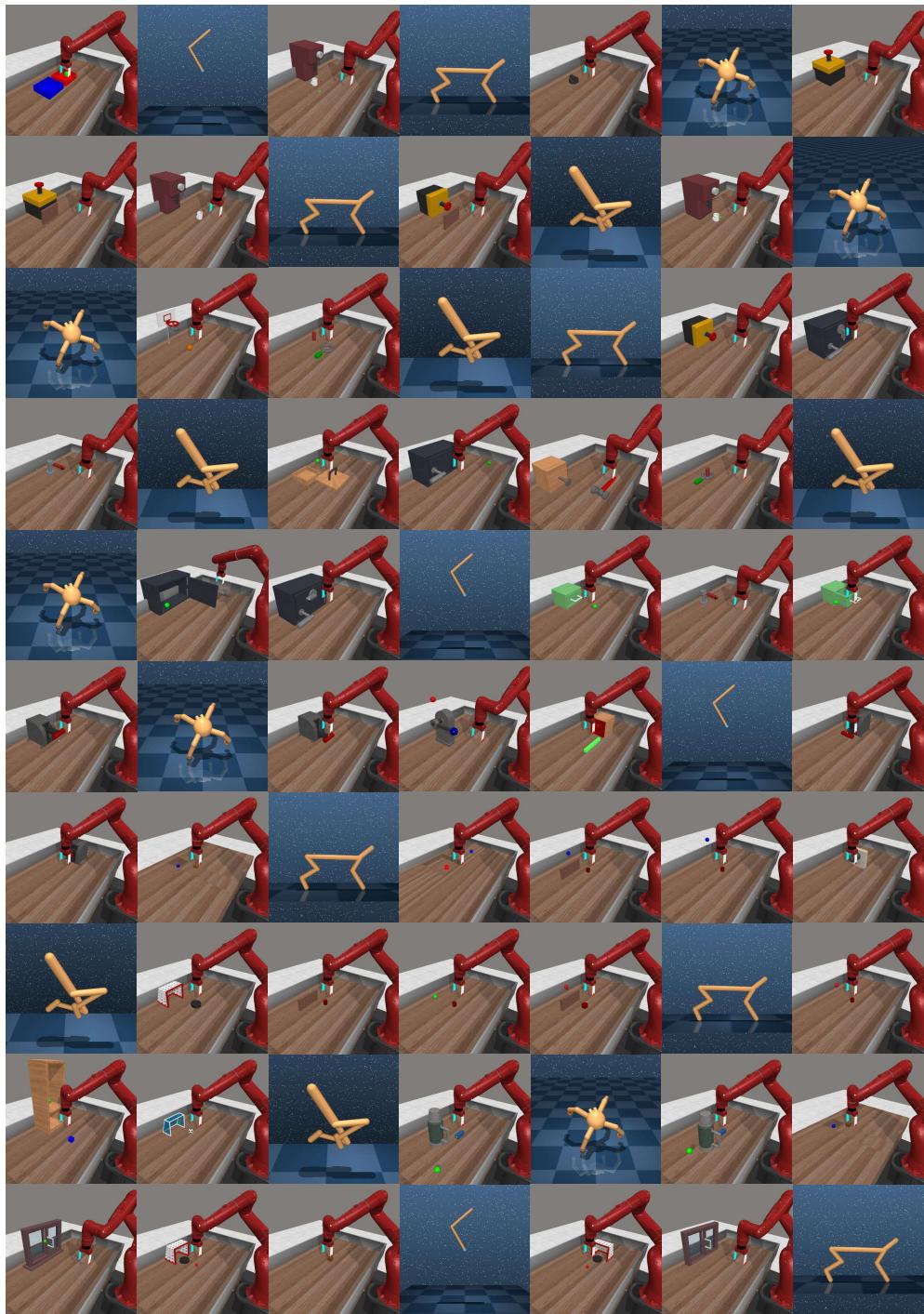


Figure 12: Visualization of tasks from DMControl and Meta-World used in our paper.