

MAKE THE PERTINENT SALIENT: TASK-RELEVANT RECONSTRUCTION FOR VISUAL CONTROL WITH DISTRACTIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent advancements in Model-Based Reinforcement Learning (MBRL) have made it a powerful tool for visual control tasks. Despite improved data efficiency, it remains challenging to train MBRL agents with generalizable perception. Training in the presence of visual distractions is particularly difficult due to the high variation they introduce to representation learning. Building on DREAMER, a popular MBRL method, we propose a simple yet effective auxiliary task to facilitate representation learning in distracting environments. Under the assumption that task-relevant components of image observations are straightforward to identify with prior knowledge in a given task, we use a segmentation mask on image observations to only reconstruct task-relevant components. In doing so, we greatly reduce the complexity of representation learning by removing the need to encode task-irrelevant objects in the latent representation. Our method, Segmentation Dreamer (SD), can be used either with ground-truth masks easily accessible in simulation or by leveraging potentially imperfect segmentation foundation models. The latter is further improved by selectively applying the reconstruction loss to avoid providing misleading learning signals due to mask prediction errors. In modified DeepMind Control suite (DMC) and Meta-World tasks with added visual distractions, SD achieves significantly better sample efficiency and greater final performance than prior work. We find that SD is especially helpful in sparse reward tasks otherwise unsolvable by prior work, enabling the training of visually robust agents without the need for extensive reward engineering.

1 INTRODUCTION

Recently, model-based reinforcement learning (MBRL) (Sutton, 1991; Ha & Schmidhuber, 2018; Hafner et al., 2019; 2020; Hansen et al., 2022; 2023) has shown great promise in learning control policies, achieving high sample efficiency. Among recent advances, the DREAMER family (Hafner et al., 2020; 2021; 2023) is considered a seminal work, showing strong performance in diverse visual control environments. This is accomplished by a close cooperation between a world model and an actor-critic agent. The world model learns to emulate the environment’s forward dynamics and reward function in a latent state space, and the agent is trained by interacting with this world model in place of the original environment.

Under this framework, accurate reward prediction is all we should sufficiently require for agent training. However, learning representations solely from reward signals is known to be challenging due to their lack of expressiveness and high variance (Hafner et al., 2020; Jaderberg et al., 2017). To address this, DREAMER employs image reconstruction as an auxiliary task in world model training to facilitate representation learning. In environments with little distraction, image reconstruction works effectively by delivering rich feature-learning signals derived from pixels. However, in the presence of distractions, the image reconstruction task pushes the encoder to retain all image information, regardless of its task relevance. For instance, moving backgrounds in observations in Fig. 1 are considered distractions. Including this information in the latent space complicates dynamics modeling and degrades sample efficiency by wasting model capacity and drowning the relevant signal in noise (Fu et al., 2021).



054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Figure 1: *Left*: Providing mask example(s) and fine-tuning a mask model, or instrumenting a simulator, to obtain masks. *Right*: An input observation in a distracting Meta-World with three alternative auxiliary task targets. Moving scenes in the background are considered distractions. (b) Observations including task-irrelevant information, disturbing world-model training. (c) and (d) Segmentation of task-relevant components using, respectively, a ground-truth mask and an approximate mask generated by segmentation models.

Distractions are prevalent in real-world visual control tasks. A robot operating in a cluttered environment such as a warehouse may perceive much task-irrelevant information that it needs to ignore. When training with domain randomization for added policy robustness, task-irrelevant information is actively added and must be denoised. Prior approaches (Zhang et al., 2021; Nguyen et al., 2021; Deng et al., 2022; Fu et al., 2021; Bharadhwaj et al., 2022) address the noisy reconstruction problem by devising reconstruction-free auxiliary tasks, such as contrastive learning (Chen et al., 2020). However, many of them suffer from sample inefficiency, requiring many trajectories to isolate the task-relevant information that needs to be encoded. Training with such methods becomes particularly challenging in sparse reward environments, where the signal for task relevance is very weak. Additionally, working with small objects, which is common in object manipulation tasks, poses difficulties for these methods because those objects contribute less to loss functions and are easily overlooked without special attention (Seo et al., 2022).

Inspired by these problems, we address the following question in this paper: *How can we help world models learn task-relevant representations more efficiently?* Our proposed solution takes advantage of the observation that identifying task-relevant components within images is often straightforward with some domain knowledge. For instance, in a robotic manipulation task, the objects to manipulate and the robot arm are such task-relevant components, as shown in Fig. 1 (Left). Given this assumption, we introduce a simple yet effective alternative auxiliary task to reconstruct only the task-related components of image observations. We accomplish this by using segmentation masks of task-related objects which are easily accessible in simulations. Specifically, we replace Dreamer’s auxiliary task to reconstruct raw RGB image observations with an alternative task to reconstruct images with a *task-relevance mask* applied to them. (Fig. 1c). By doing this, the world model can learn features from a rich pixel-reconstruction loss signal without being hindered by the noise of visual distractions.

In contrast to previous work that incorporates segmentation masks in reinforcement learning (RL) as an input (James et al., 2019; So et al., 2022), we only use masks in an auxiliary task to improve representation learning. This brings about several advantages. First, we only need segmentation masks during training. The inputs for our method are still the original (potentially distracting) images, so masks are not needed at test time, making our method more computationally efficient. Moreover, the masks we use do not need to be perfect. Since they are used as a target for an auxiliary task, these masks can contain errors as long as they guide feature learning to be informative for the downstream task, which leaves room to replace a ground-truth (GT) mask with its approximation.

To this end, we present a way of training with our auxiliary task with segmentation estimates. This can be useful in many practical cases where no GT mask is available during training. This is made possible by recent advances in segmentation foundation models (Kirillov et al., 2023; Zhang et al., 2023; Xie et al., 2021). Specifically, we fine-tune segmentation models with annotated training data and use them to generate pseudo-labels for the auxiliary task. Fig. 1d shows an example of an auxiliary target made from segmentation prediction. Although the performance with segmentation estimates is impressive without further modification, we find that the training can sometimes

108 be destabilized by incorrect learning signals induced by segmentation prediction errors. Thus, we
109 additionally provide a strategy to make training more robust to prediction errors and achieve higher
110 performance. Specifically, our strategy is to identify pixels where the foundation model mask pre-
111 diction disagrees with a second mask prediction given by our world model. We ignore the RGB
112 image reconstruction loss on these pixels to avoid training on potentially incorrect targets.

113 As previously mentioned, our method assumes that task-relevant parts are easy to identify in im-
114 age observations with prior knowledge. This is not a strong assumption in many object-centric and
115 robotic domains, where image observations can often be decomposed into relevant and irrelevant re-
116 gions. However, there are scenarios beyond our scope, where this assumption may not hold because
117 prior knowledge is unavailable or difficult to codify, such as Atari (Bellemare et al., 2013).

118 We evaluate our method on various robotics benchmarks, including DeepMind Control Suite (Tassa
119 et al., 2018) and Meta-World (Yu et al., 2019), perturbing both with visual distractions. We show
120 that our method for reconstructing masked RGB targets using the ground-truth masks in the presence
121 of distractions can reach the same level of performance as training in *original* environment with no
122 distractions added. Our method for training with approximate masks also shows impressive perfor-
123 mance, often matching the performance of the ground-truth mask variant. In both benchmarks, our
124 approximate-mask method achieves higher sample efficiency and superior test returns compared to
125 previous approaches. Notably, this is accomplished with very few task-specific mask example data
126 points (1, 5, or 10 used for fine-tuning), with much of its strength coming from the power of seg-
127 mentation foundation models. Our method effectively addresses the challenge of training agents in
128 distracting environments by offloading the identification of task-relevant regions to out-of-the-box
129 segmentation models, thereby achieving great sample efficiency and generalization ability. Further-
130 more, our method proves particularly effective in sparse reward environments and those involving
131 small objects, where prior approaches often struggle.

132 2 RELATED WORK

135 **Model-Based RL for Distracting Visual Environments.** Recent advances in MBRL have facil-
136 itated the learning of agents from image observations (Finn & Levine, 2017; Ha & Schmidhuber,
137 2018; Hafner et al., 2019; 2020; 2021; 2023; Schrittwieser et al., 2020; Hansen et al., 2022; 2023).
138 Nevertheless, learning perceptual representations in the presence of distractions remains challeng-
139 ing in these models (Zhu et al., 2023). For effective representation learning, some works apply
140 non-reconstructive representation learning methods (Nguyen et al., 2021; Deng et al., 2022), such as
141 contrastive learning (Chen et al., 2020) and prototypical representation learning (Caron et al., 2020).
142 However, features learned with these methods do not necessarily involve task-related content since
143 they do not explicitly consider task-relevance in feature learning. Some other works design auxil-
144 iary objectives to explicitly use downstream task information (Zhang et al., 2021; Fu et al., 2021).
145 For example, DBC (Zhang et al., 2021) uses a bisimulation metric (Ferns et al., 2011) to encourage
146 two trajectories of similar behaviors to become closer in a latent space. Perhaps most relevant to
147 our method is TIA (Fu et al., 2021) which explicitly separates task-relevant and irrelevant branches
148 to distinguish reward-correlated visual features from distractions. Features from each branch are
149 combined later to reconstruct the original, distracting observation. Recently, a few approaches pro-
150 posed to leverage inductive biases such as predictability (Zhu et al., 2023) and controllability (Wang
151 et al., 2022; Bharadhwaj et al., 2022) to learn useful features for visual control tasks. These methods
152 have shown to be more effective than using the reward signal alone, but many of them suffer from
153 sample inefficiency, requiring many samples to implicitly identify what is task-relevant from data.
154 In contrast, our work proposes to leverage domain knowledge in the form of image masks to pro-
155 vide an explicit signal for identifying task-relevant information. Notably, training in sparse reward
156 environments with distraction has remained unsolved in the literature. Several methods for robust
157 representation learning have also been proposed for model-free RL (Laskin et al., 2020; Kostrikov
158 et al., 2021; Yarats et al., 2021; Hansen et al., 2021; Hansen & Wang, 2021; Nair et al., 2022; Zhang
159 et al., 2019). However, the results suggest that MBRL is more powerful and sample efficient for
160 visual control tasks, thus we focus on comparison with methods in the MBRL framework.

161 **Segmentation for RL.** Segmentation models (He et al., 2017; Redmon et al., 2016) have been
used in many downstream tasks, including RL, to assist in pre-processing inputs (Kirillov et al.,
2023; Anantharaman et al., 2018; Yuan et al., 2018; James et al., 2019; So et al., 2022). Recently,

using segmentation masks in new domains has been made easier by the introduction of one/few-shot masks foundation models (Zhang et al., 2023; Xie et al., 2021) which can quickly adapt to new use-cases. In the context of RL, the prevalent way to use segmentation models is to turn the input modality from RGB images to segmentation masks (James et al., 2019; So et al., 2022). By converting RGB images into semantic masks, agents can effectively handle complicated scenes and thus also be trained with domain randomization. However, processing the input observation requires additional computation at execution time, and an agent trained with predicted segmentation masks can easily break when the segmentation model malfunctions. Our approach, on the other hand, leverages segmentation masks for an auxiliary task, removing the need for the segmentation model after deployment while also achieving higher test-time performance. FOCUS (Ferraro et al., 2023) is another method that uses masked input as an auxiliary target. However, this is primarily devised for disentangled representation, not for handling distractions. Moreover, it only provides preliminary results with segmentation models and lacks results and analysis on the effects on downstream tasks.

Prior works (Wang et al., 2023; Zhong et al., 2024) integrate segmentation models with RL by pre-processing input observations to isolate task-relevant components. While effective, these methods heavily rely on segmentation model quality at test time, making them vulnerable to failures in unfamiliar scenarios that disrupt agent performance. Also, methods such as Zhong et al. (2024); So et al. (2022) require extensive fine-tuning on large synthetic datasets, resulting in substantial training overhead, and introduce high computational costs during inference. By leveraging off-the-shelf segmentation models with as few as 1 to 10 examples, our approach reduces training requirements while maintaining robustness and runtime efficiency during deployment.

3 PRELIMINARIES

We consider a partially observable Markov decision process (POMDP) formalized as a tuple $(\mathcal{S}, \Omega, \mathcal{A}, \mathcal{T}, \mathcal{O}, p_0, \mathcal{R}, \gamma)$, consisting of states $s \in \mathcal{S}$, observations $o \in \Omega$, actions $a \in \mathcal{A}$, state transition function $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$, observation function $\mathcal{O} : \mathcal{S} \rightarrow \Omega$, initial state distribution p_0 , reward function $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$, and discount factor γ . At time t , the agent does not have access to actual world state s_t , but to the observation $o_t = \mathcal{O}(s_t)$, which in this paper we consider to be a high-dimensional image. Our objective is to learn a policy $\pi(a_t | o_{\leq t}, a_{< t})$ that achieves high expected discounted cumulative rewards $\mathbb{E}[\sum_t \gamma^t r_t]$, with $r_t = \mathcal{R}(s_t, a_t)$ and the expectation over the joint stochastic process induced by the environment and the policy.

DREAMER (Hafner et al., 2020; 2021; 2023) is a broadly applicable MBRL method in which a world model is learned to represent environment dynamics in a latent state space $(h, z) \in \mathcal{H} \times \mathcal{Z}$, consisting of deterministic and stochastic components respectively, from which rewards, observations, and future latent states can be decoded. The components of the world model are:

$$\begin{aligned}
 \text{Sequence model:} & \quad h_t = f_\phi(h_{t-1}, z_{t-1}, a_{t-1}) \\
 \text{Observation encoder:} & \quad z_t \sim q_\phi(z_t | h_t, o_t) \\
 \text{Dynamics predictor:} & \quad \hat{z}_t \sim p_\phi(\hat{z}_t | h_t) \\
 \text{Reward predictor:} & \quad \hat{r}_t \sim p_\phi(\hat{r}_t | h_t, z_t) \\
 \text{Continuation predictor:} & \quad \hat{c}_t \sim p_\phi(\hat{c}_t | h_t, z_t) \\
 \text{Observation decoder:} & \quad \hat{o}_t \sim p_\phi(\hat{o}_t | h_t, z_t),
 \end{aligned} \tag{1}$$

where the encoder maps observations o_t into a latent representation, the dynamics model emulates the transition distribution in latent state space, the reward and continuation models respectively predict rewards and episode termination, and the observation decoder reconstructs the input. The concatenation of h_t and z_t , *i.e.* $x_t = [h_t; z_t]$, serves as the model state. Given a starting state, an actor-critic agent is trained inside the world model by rolling out latent-state trajectories. The world model itself is trained by optimizing a weighted combination of three losses:

$$\mathcal{L}(\phi) \doteq \mathbb{E}_{q_\phi} \left[\sum_{t=1}^T (\beta_{\text{pred}} \mathcal{L}_{\text{pred}}(\phi) + \beta_{\text{dyn}} \mathcal{L}_{\text{dyn}}(\phi) + \beta_{\text{rep}} \mathcal{L}_{\text{rep}}(\phi)) \right] \tag{2}$$

$$\mathcal{L}_{\text{pred}}(\phi) \doteq -\ln p_\phi(o_t | z_t, h_t) - \ln p_\phi(r_t | z_t, h_t) - \ln p_\phi(c_t | z_t, h_t) \tag{3}$$

$$\mathcal{L}_{\text{dyn}}(\phi) \doteq \max(1, \text{KL}[[q_\phi(z_t | h_t, o_t)] \| p_\phi(\hat{z}_t | h_t)]) \tag{4}$$

$$\mathcal{L}_{\text{rep}}(\phi) \doteq \max(1, \text{KL}[q_\phi(z_t | h_t, o_t) \| [p_\phi(\hat{z}_t | h_t)]]) \tag{5}$$

where $\llbracket \cdot \rrbracket$ denotes where gradients are stopped from backpropagating to the expression in brackets.

Critically, the first component of $\mathcal{L}_{\text{pred}}$ for reconstructing observations from world model states is leveraged as a powerful heuristic to shape the features in the latent space. Under the assumption that observations primarily contain task-relevant information, this objective is likely to encourage the latent state to retain information critical for the RL agent. However, the opposite can also be true. If observations are dominated by task-irrelevant information, the latent dynamics may become more complex by incorporating features impertinent to decision-making. This can lead to wasted capacity in the latent state representation (Lambert et al., 2020), drown the supervision signal in noise, and reduce the sample efficiency.

Problem Setup. We consider environments where the latter case is true and observations contain a large number of spurious variations (Zhu et al., 2023). Concretely, we consider some features of states $s_t \in \mathcal{S}$ to be irrelevant for the control task. We assume that states s_t can be decomposed into task-relevant components $s_t^+ \in S^+$ and task-irrelevant components $s_t^- \in S^-$ such that $s_t = (s_t^+, s_t^-) \in \mathcal{S} = S^+ \times S^-$. We follow prior work in visual control under distraction and assume that (1) the reward is a function only of the task-relevant component, *i.e.* $\mathcal{R} : S^+ \times \mathcal{A} \rightarrow \mathbb{R}$; and (2) the forward dynamics of the task-relevant part only depends on itself, $s_{t+1}^+ \sim \mathcal{T}(s_{t+1}^+ | s_t^+, a_t)$ (Zhu et al., 2023; Fu et al., 2021; Bharadhwaj et al., 2022). Note that observations o_t are a function of both s_t^+ and s_t^- , thus we have $\mathcal{O} : S^+ \times S^- \rightarrow \Omega$.

Our goal is to learn effective latent representations $[h_t; z_t]$ for task control. Ideally, this would mean that the world model will only encode and simulate task-relevant state components s_t^+ in its latent space without modeling unnecessary information in s_t^- . To learn features pertaining to s_t^+ , image reconstruction can provide a rich and direct learning signal, but only when observation information about s_t^+ is not drowned out by other information from s_t^- . To overcome this pitfall, we propose to apply a heuristic filter to reconstruction targets o_t with the criteria that it minimizes irrelevant information pertaining to s_t^- while keeping task-relevant information about s_t^+ .

4 METHOD

We build on DREAMER-V3 (Hafner et al., 2023) to explicitly model s_t^+ while attempting to avoid encoding information about s_t^- . In Section 4.1, we describe how we accomplish this by using domain knowledge to apply a task-relevance mask to observation reconstruction targets. In Section 4.2 we describe how we leverage segmentation mask foundation models to provide approximate masks over task-relevant observation components. Finally, in Section 4.3, we propose a modified decoder architecture and objective to mitigate noisy learning signals from incorrect mask predictions.

4.1 USING SEGMENTATION MASKS TO FILTER IMAGE TARGETS

We first introduce our main assumption, that the task-relevant components of image observations are easily identifiable with domain knowledge. In many real scenarios, it is often straightforward for a practitioner to know what the task-related parts of an image are, *e.g.* objects necessary for achieving a goal in object manipulation tasks. With this assumption, we propose a new reconstruction-based auxiliary task that leverages domain knowledge of task-relevant regions. Instead of reconstructing the raw image observations (Fig. 1b) which may contain task-irrelevant distractions, we apply a heuristic task-relevance segmentation mask over the image observation (Fig. 1c) to exclusively reconstruct components of the image that are pertinent to control.

Since our new masked reconstruction target should contain only image regions that are relevant for achieving the downstream task, our world model should learn latent representations where a larger portion of the features are useful to the RL agent. By explicitly avoiding modeling task-irrelevant observation components, the latent dynamics should also become simpler and more sample-efficient to learn than the original (more complex, higher variance) dynamics on unfiltered observations. In simulations, ground-truth masks of relevant observation components are often easily accessible, for example, in MuJoCo (Todorov et al., 2012), through added calls to the simulator API. We term the method trained with our proposed replacement auxiliary task as Segmentation Dreamer (SD) and call the version trained with ground-truth masks SD^{GT} .

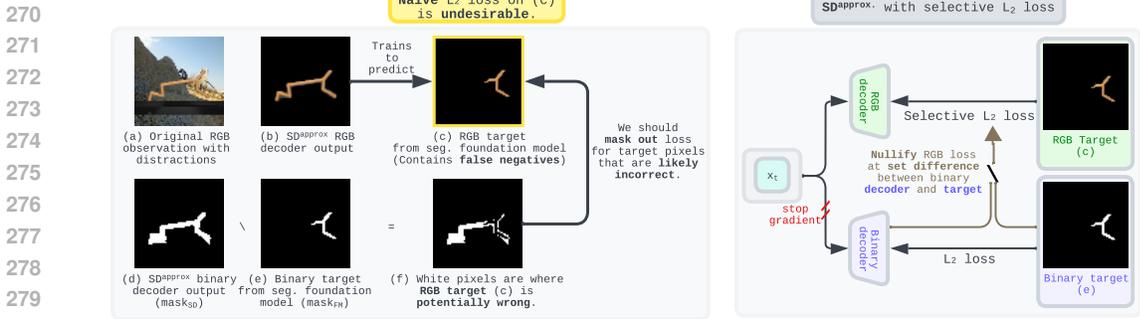


Figure 2: **Filtering L_2 loss to avoid training on false negatives in RGB labels.** *Left:* Estimated pixel locations (f) where the RGB target (c) is likely incorrectly masked out by the segmentation model (e). *Right:* A world model equipped with two decoders, one for reconstructing task-relevant masked RGB images and the other for binary masks, the targets for which are generated by a segmentation model. RGB L_2 loss is selectively masked by the set difference between (d) and (e). Latent representations (x_t) in the world model are subjected to the training signal only from the RGB branch. The binary branch is only utilized for selective L_2 loss.

4.2 LEVERAGING APPROXIMATE SEGMENTATION MASKS

A simulator capable of providing ground-truth masks for task-relevant regions is not always available. For such cases where only RGB images are available from the environment, we propose to fine-tune a segmentation mask foundation model to our domain and integrate its predictions into the SD training pipeline. Below, we describe our method for training with approximate task-relevance masks, termed SD^{approx} .

As an offline process before training the world model, we fine-tune a segmentation model with a small number of example RGB images and their segmentation masks annotations that indicate task-relevant image regions. Thanks to recent advances in segmentation foundation models, we can obtain a new domain-specific mask model with a very small amount of training examples. For our experiments, we use the Personalized SAM (PerSAM) (Zhang et al., 2023) using one-shot adaptation and SegFormer (Xie et al., 2021) fine-tuned with 5 and 10 examples. For the sake of controlled and reproducible evaluation, we extract these RGB and mask training pairs from simulators, however, this number of required samples is small enough that it can be collected with expert human annotation as well. Also, although we use these specific foundation models, our method should also be compatible with *any* semantic masking method. Further details such as how we obtain the fine-tuning data can be found in Appendix K. Once fine-tuning is complete, we incorporate the segmentation model into the SD pipeline to create pseudo-labels for our proposed auxiliary task.

4.3 LEARNING IN THE PRESENCE OF MASKING ERRORS

Although foundation segmentation models generalize well to new scenarios (e.g., different poses, occlusions), prediction errors are inevitable (Fig. 1d). Since each frame is processed independently, segmentation predictions can flicker along trajectories. False negatives in task relevance are particularly detrimental when using naive L_2 loss on image reconstruction. Missing relevant scene elements in reconstruction targets can lead the encoder to learn incomplete representations, dropping essential task-related information. This variability disrupts the learning of accurate representations and dynamics in the world model.

Despite noisy targets, neural networks can self-correct if most labels are accurate (Han et al., 2018). Additionally, DREAMER’s use of GRUs (Cho et al., 2014) provides temporal consistency even with flickering targets. However, as illustrated in Fig. 2 (b)&(c), it’s undesirable to propagate gradients from regions where the original image has been incorrectly masked out. Allowing gradients from these regions provides misleading signals. If we could identify the incorrect regions in the reconstruction target, we could nullify the decoder’s L_2 loss there—a technique we call selective L_2 loss.

Since we cannot directly identify regions where the RGB target is incorrectly masked due to false negatives, we estimate them. Preliminary experiments show that a binary mask decoder from world

model states (as an added auxiliary task) can be less prone to transient false negatives, unlike RGB prediction, which tends to memorize noisy labels. Therefore, we propose training a world model with two reconstruction tasks (Fig. 2, right): one decoding masked RGB images and the other predicting task-relevance binary masks. Both use the foundation model’s binary mask, mask_{FM} , to construct targets. The RGB branch decodes masked RGB images, while the binary branch predicts mask_{FM} . We denote the binary masks produced by the world model as mask_{SD} , where pixels labeled *true* indicate task relevance.

To avoid training on incorrectly masked-out regions, we estimate where mask_{FM} may be falsely negative by finding disagreements with mask_{SD} . Specifically, we selectively nullify RGB decoder L_2 loss for regions marked false in mask_{FM} but predicted true in mask_{SD} . This prevents training on potentially falsely masked-out pixels still considered task-relevant by a second predictor. Formally, the mask for selective L_2 loss is the set difference between true pixel locations in mask_{SD} and mask_{FM} :

$$\text{pixel}_{\text{MaskOut}} = \text{pixel}_{\text{SD}} \setminus \text{pixel}_{\text{FM}} \quad (6)$$

where $\text{pixel}_{\text{MaskOut}}$ indicates pixels to nullify loss at, and pixel_{SD} and pixel_{FM} are pixels marked true in mask_{SD} and mask_{FM} , respectively.

Fig. 2 (d–f) shows examples of mask_{SD} , mask_{FM} , and $\text{pixel}_{\text{MaskOut}}$. See Appendix L for details on obtaining mask_{SD} . Our experiments indicate that selective L_2 loss effectively overcomes noisy segmentation labels and improves downstream agent performance.

Lastly, we observe better performance when we prevent gradients from the binary mask decoding objective from propagating into the world model, so we apply a stop gradient to the inputs of the mask decoder head (see Appendix G for ablations).

5 EXPERIMENTS

We evaluate our method on a variety of visual robotic control tasks from the DeepMind Control Suite (DMC) (Tassa et al., 2018) and Meta-World (Yu et al., 2019). Since the standard environments in these benchmarks have simple backgrounds with minimal distractions, we introduce visual distractions by replacing the backgrounds with random videos from the ‘driving car’ class in the Kinetics 400 dataset (Kay et al., 2017), following prior work (Zhang et al., 2021; Nguyen et al., 2021; Deng et al., 2022). Details about the environment setup and task visualizations are provided in Appendices H and B. In evaluation, we roll out policies over 10 episodes and compute the average episode return. Unless otherwise specified, we report the mean and standard error of the mean (SEM) of four independent runs with different random seeds. We use default DREAMER-V3 hyperparameters in all experiments.

5.1 DMC EXPERIMENTS

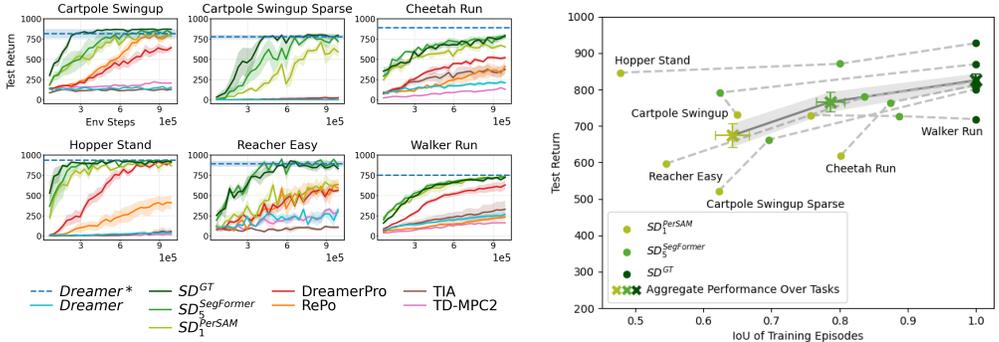
We evaluate SD on six tasks from DMC featuring different forms of contact dynamics, degrees of freedom, and reward sparsities. For each task, models are trained for 1M environment steps generated by 500K policy decision steps with an action repeat of 2.

5.1.1 COMPARISON WITH DREAMER

We compare our methods, SD^{GT} and $\text{SD}^{\text{approx.}}$, to the base DREAMER (Hafner et al., 2023) method. Here, $\text{SD}^{\text{approx.}}$ is denoted as SD_N^{FM} , specifying the segmentation model used (FM) and the number of fine-tuning examples (N). All methods are trained in distracting environments, except for the DREAMER* baseline, which is trained in the original environment without visual distractions. In most cases, we consider DREAMER* as an upper bound for methods trained with distractions. Similarly, SD^{GT} serves as an upper bound for $\text{SD}^{\text{approx.}}$, with the performance gap expected to decrease in the future as segmentation quality improves.

As shown in Fig. 3a, DREAMER fails across all tasks due to task-irrelevant information in RGB reconstruction targets, which wastes latent capacity and complicates dynamics learning. In contrast, SD^{GT} achieves test returns comparable to DREAMER* by focusing on reconstructing essential features and ignoring irrelevant components. Interestingly, SD^{GT} outperforms DREAMER* in Cartpole

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431



(a) Environment Steps vs. Expected Test Return (b) IoU during Training vs. Expected Test Return

Figure 3: (a) **Learning curves on six visual control tasks from DMC.** Every method but DREAMER* is trained on distracting environments. All curves show the mean over 4 seeds with the standard error of the mean (SEM) shaded. (b) Segmentation quality during training vs. downstream task performance. Best viewed in color.

Swingup, possibly because the original environment still contains small distractions (e.g., moving dots) that DREAMER* has to model.

A limitation of SD is its reliance on accurate and correct prior knowledge to select task-relevant components. In Cheetah Run, SD^{GT} underperforms compared to DREAMER*, likely because we only include the cheetah’s body in the mask, excluding the ground plate, which may be important for contact dynamics. Visual examples and further experiments are in Appendices B and C.

For SD^{approx.}, we test with two foundation models: PerSAM adapted with one RGB example and its GT mask, and SegFormer adapted with five such examples. Despite slower convergence due to noisier targets, both SD₁^{PerSAM} and SD₅^{SegFormer} achieve similar final performance to SD^{GT} in most tasks. A failure case for SD₁^{PerSAM} is Reacher Easy, where a single data point is insufficient to obtain a quality segmentation for the small task-relevant objects.

5.1.2 COMPARISON WITH BASELINES

We compare SD^{approx.} with state-of-the-art methods, including DreamerPro (Deng et al., 2022), RePo (Zhu et al., 2023), TIA (Fu et al., 2021), and TD-MPC2 (Hansen et al., 2023). DreamerPro incorporates prototypical representation learning in the DREAMER framework; RePo minimizes mutual information between observations and latent states while maximizing it between states and future rewards; TIA learns separate task-relevant and task-irrelevant representations which can be combined to decode observations; and TD-MPC2 decodes a terminal value function. Among these baselines, only TIA relies on observation reconstruction. Further details are in Appendix M.

Our results in Fig. 3a show that our method consistently outperforms the baselines in performance and sample efficiency. TIA underperforms in many tasks, requiring many samples to infer task-relevant observations from rewards and needing exhaustive hyperparameter tuning. Even with optimal settings, it may lead to degenerate solutions where a single branch captures all information. In contrast, our method focuses on task-relevant parts without additional tuning by effectively injecting prior knowledge. RePo performs comparably to ours in Cartpole Swingup but underperforms in other tasks and converges more slowly.

TD-MPC2 struggles significantly in distracting environments. We speculate that spurious correlations from distractions introduce noise to value-function credit assignment that hinders representation learning. Our method mitigates this by directly supervising task-relevant features using segmentation models, leading to more consistent and lower-variance targets.

Among these methods, DreamerPro is the most competitive, demonstrating the effectiveness of prototypical representation learning for control. However, it often requires more environment interactions and converges to lower performance.

In the Cartpole Swingup with sparse rewards, none of the prior works successfully solved the task, highlighting the challenge of inferring task relevance from weak signals. Our method achieves near-

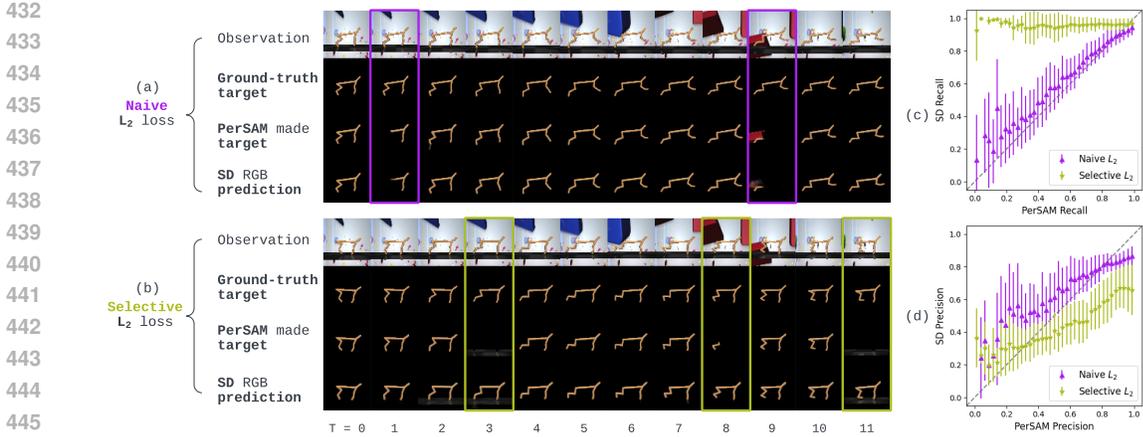


Figure 4: (a)+(b) **Qualitative comparison of SD trained with naive and selective L_2 loss.** Trajectories are taken from each method’s train-time replay buffer, selected to have the same background. Frames with PerSAM error are highlighted. The model trained with the selective L_2 loss overcomes errors in the target, whereas the one trained with the naive L_2 loss memorizes target errors. (c)+(d) shows the precision and recall of PerSAM and the SD RGB decoder prediction. SD RGB predictions are binarized using a threshold to compute recall and precision w.r.t. the ground-truth mask. The data points used for plotting are from the same Cheetah Run training experiment as in (a)+(b). The selective L_2 loss significantly improves the recall with only a moderate impact on precision.

oracle performance, being the only one to train an agent with sparse rewards amidst distractions. This suggests the potential to train agents in real-world, distraction-rich environments without extensive reward engineering.

5.1.3 ABLATION STUDY

We investigate the effects of the components in SD^{approx} . by addressing: (1) the benefits of using segmentation models for targets vs. input preprocessing; (2) the effectiveness of the selective L_2 loss compared to the naive L_2 loss; and (3) the impact of the segmentation quality on RL performance. In these experiments, we fine-tune PerSAM with a single data point for segmentation mask prediction.

Using segmentation masks for an auxiliary task vs. input preprocessing. We create a variant of SD_1^{PerSAM} that uses masked observations for both inputs and targets, denoted in Tab. 1 by *As Input*. These results suggest that SD_1^{PerSAM} , in addition to not requiring mask prediction at test-time, also achieves better test performance and lower variance. Using predicted masks as input is more prone to segmentation errors, restricting the agent’s perception when masks are incorrect and making training more challenging. In contrast, SD^{approx} . receives intact observations, with task-relevant filtering at the encoder level, leading to better state abstraction. Further analysis on test-time segmentation quality’s impact is in Appendix D.

Table 1: Final performance of SD variants. Mean over 4 runs with the standard error of the mean is reported. The highest means are highlighted.

Task	SD_1^{PerSAM}	As Input	Naive L_2
Cartpole Swingup	730 ± 75	565 ± 108	719 ± 62
Cartpole Swingup Sparse	521 ± 92	457 ± 151	408 ± 114
Cheetah Run	619 ± 35	524 ± 37	486 ± 58
Hopper Stand	846 ± 27	689 ± 39	790 ± 51
Reacher Easy	597 ± 97	642 ± 116	415 ± 50
Walker Run	730 ± 13	589 ± 28	557 ± 51

Selective L_2 loss vs. naive L_2 loss. As shown in Tab. 1, SD_1^{PerSAM} consistently outperforms the Naive L_2 variant, especially in complex tasks like Cheetah Run and Walker Run. Segmentation models often miss embodiment components (Fig. 4, third row). With the naive L_2 loss, the model replicates these errors, leading to incomplete latent representations and harming dynamics learning (Fig. 4a, fourth row). In contrast, SD^{approx} . self-corrects by skipping the L_2 computation where PerSAM targets are likely wrong (Fig. 4b, fourth row). Fig. 4(c)&(d) show that the naive L_2 loss

follows PerSAM’s trends, while the selective L_2 loss recovers from poor recall with only a moderate precision decrease.

Impact of segmentation quality on RL performance. Fig. 3b plots the training-time segmentation quality against the RL agent’s test-time performance. Comparing three SD variants with different mask qualities (two estimated, one ground truth), we observe that better segmentation tends to lead to higher RL performance, as accurate targets better highlight task-relevant components. This suggests that improved segmentation models can enhance agent performance without ground-truth masks. In Cartpole Swingup, one of two exceptions, the IoU difference between SD_1^{PerSAM} and $SD_5^{\text{SegFormer}}$ is small, and the test returns may fall within the margin of error. In Walker Run, the other exception, all variants show high segmentation quality and reach near-optimal performance. Here, we hypothesize that a small amount of noise in the target may act as a regularizer, contributing to marginally better downstream performance.

5.2 META-WORLD EXPERIMENTS

Object manipulation is a natural application for our method where prior knowledge can be applied straightforwardly by identifying and masking task-relevant objects and robot embodiments. We evaluate SD on six tasks from Meta-World (Yu et al., 2019), a popular benchmark for robotic manipulation. Depending on the difficulty of each task, we conduct experiments for 30K, 100K, and 1M environment steps, with an action repeat of 2 (details in Appendix I). Preliminary tests showed that SegFormer performs well with few-shot learning on small objects. We fine-tune SegFormer with 10 data points to estimate masks in these experiments.

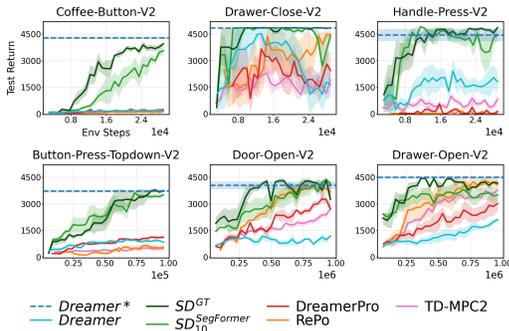


Figure 5: **Learning curves on six visual robotic manipulation tasks from Meta-World.** All curves show the mean over 4 seeds with the standard error of the mean shaded.

Fig. 5 suggests that our approach outperforms the baselines overall, with a more pronounced advantage in tasks involving small objects like Coffee-Button. Our method excels because it focuses on small, task-relevant objects, avoiding the reconstruction of unnecessary regions that occupy much of the input. In contrast, the baselines struggle as they often underestimate the significance of these small yet highly task-relevant objects. Among the baselines, RePo (Zhu et al., 2023) is the most competitive. However, RePo performs poorly in a sparse reward setup (see Appendix J).

6 CONCLUSION

In this paper, we propose SD, a simple yet effective method for learning task-relevant features in MBRL frameworks like DREAMER by using segmentation masks informed by domain knowledge. Using ground-truth masks, SD^{GT} achieves performance comparable with undistracted DREAMER with high sample efficiency in distracting environments when provided with accurate prior knowledge. Our main method, SD^{approx} , uses mask estimates from off-the-shelf one-shot or few-shot segmentation models and employs a selective L_2 loss. It learns effective world models that produce strong agents outperforming baselines.

To the best of our knowledge, our approach appears to be the first model-based approach to successfully train an agent in a sparse reward environment under visual distractions, enabling robust agent training without extensive reward engineering. This work also advances the integration of computer vision and RL by presenting a novel way to leverage recent advances in segmentation to address challenges in visual control tasks. The proposed method achieves strong performance on diverse tasks with distractions and effectively incorporates human input to indicate task relevance. This enables practitioners to readily train an agent for their own purposes without extensive reward engineering. However, SD has some limitations to consider in future work, which we further explore in Appendix O.

REFERENCES

- 540
541
542 Rajaram Anantharaman, Matthew Velazquez, and Yugyung Lee. Utilizing mask r-cnn for de-
543 tection and segmentation of oral diseases. In *International Conference on Bioinformatics and*
544 *Biomedicine*, pp. 2197–2204. IEEE, 2018.
- 545 Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environ-
546 nment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:
547 253–279, 2013.
- 548 Homanga Bharadhwaj, Mohammad Babaeizadeh, Dumitru Erhan, and Sergey Levine. Information
549 prioritization through empowerment in visual model-based rl. In *International Conference on*
550 *Learning Representations*, 2022.
- 551 Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin.
552 Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural*
553 *Information Processing Systems*, 33:9912–9924, 2020.
- 554 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
555 contrastive learning of visual representations. In *International Conference on Machine Learning*,
556 pp. 1597–1607. PMLR, 2020.
- 557 Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the prop-
558 erties of neural machine translation: Encoder–decoder approaches. In Dekai Wu, Marine
559 Carpuat, Xavier Carreras, and Eva Maria Vecchi (eds.), *Proceedings of SSST-8, Eighth Work-*
560 *shop on Syntax, Semantics and Structure in Statistical Translation*, pp. 103–111, Doha, Qatar,
561 October 2014. Association for Computational Linguistics. doi: 10.3115/v1/W14-4012. URL
562 <https://aclanthology.org/W14-4012>.
- 563 Fei Deng, Ingo Jang, and Sungjin Ahn. Dreamerpro: Reconstruction-free model-based rein-
564 forcement learning with prototypical representations. In *International Conference on Machine*
565 *Learning*, pp. 4956–4975. PMLR, 2022.
- 566 Norm Ferns, Prakash Panangaden, and Doina Precup. Bisimulation metrics for continuous markov
567 decision processes. *SIAM Journal on Computing*, 40(6):1662–1714, 2011.
- 568 Stefano Ferraro, Pietro Mazzaglia, Tim Verbelen, and Bart Dhoedt. Focus: Object-centric world
569 models for robotics manipulation. *arXiv preprint arXiv:2307.02427*, 2023.
- 570 Chelsea Finn and Sergey Levine. Deep visual foresight for planning robot motion. In *International*
571 *Conference on Robotics and Automation*, pp. 2786–2793. IEEE, 2017.
- 572 Xiang Fu, Ge Yang, Pulkit Agrawal, and Tommi Jaakkola. Learning task informed abstractions. In
573 *International Conference on Machine Learning*, pp. 3480–3491. PMLR, 2021.
- 574 Carles Gelada, Saurabh Kumar, Jacob Buckman, Ofir Nachum, and Marc G Bellemare. Deepmdp:
575 Learning continuous latent space models for representation learning. In *International Conference*
576 *on Machine Learning*, pp. 2170–2179. PMLR, 2019.
- 577 Jiayuan Gu, Fanbo Xiang, Xuanlin Li, Zhan Ling, Xiqiang Liu, Tongzhou Mu, Yihe Tang, Stone
578 Tao, Xinyue Wei, Yunchao Yao, et al. Maniskill2: A unified benchmark for generalizable manip-
579 ulation skills. In *International Conference on Learning Representations*, 2023.
- 580 David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- 581 Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James
582 Davidson. Learning latent dynamics for planning from pixels. In *International Conference on*
583 *Machine Learning*, pp. 2555–2565. PMLR, 2019.
- 584 Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning
585 behaviors by latent imagination. In *International Conference on Learning Representations*, 2020.
- 586 Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with dis-
587 crete world models. In *International Conference on Learning Representations*, 2021.

- 594 Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains
595 through world models. *arXiv preprint arXiv:2301.04104*, 2023.
596
- 597 Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor Tsang, and Masashi
598 Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels.
599 *Advances in Neural Information Processing Systems*, 31, 2018.
- 600 Nicklas Hansen and Xiaolong Wang. Generalization in reinforcement learning by soft data aug-
601 mentation. In *International Conference on Robotics and Automation*, pp. 13611–13617. IEEE,
602 2021.
- 603 Nicklas Hansen, Hao Su, and Xiaolong Wang. Stabilizing deep q-learning with convnets and vision
604 transformers under data augmentation. In *Advances in Neural Information Processing Systems*,
605 volume 34, pp. 3680–3693, 2021.
- 607 Nicklas Hansen, Xiaolong Wang, and Hao Su. Temporal difference learning for model predictive
608 control. In *International Conference on Machine Learning*, 2022.
609
- 610 Nicklas Hansen, Hao Su, and Xiaolong Wang. Td-mpc2: Scalable, robust world models for contin-
611 uous control. *arXiv preprint arXiv:2310.16828*, 2023.
- 612 Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *International
613 Conference on Computer Vision*, pp. 2961–2969, 2017.
614
- 615 Max Jaderberg, Volodymyr Mnih, Wojciech Marian Czarnecki, Tom Schaul, Joel Z Leibo, David
616 Silver, and Koray Kavukcuoglu. Reinforcement learning with unsupervised auxiliary tasks. In
617 *International Conference on Learning Representations*, 2017.
- 618 Stephen James, Paul Wohlhart, Mrinal Kalakrishnan, Dmitry Kalashnikov, Alex Irpan, Julian Ibarz,
619 Sergey Levine, Raia Hadsell, and Konstantinos Bousmalis. Sim-to-real via sim-to-sim: Data-
620 efficient robotic grasping via randomized-to-canonical adaptation networks. In *Computer Vision
621 and Pattern Recognition*, pp. 12627–12637, 2019.
622
- 623 Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijaya-
624 narasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action
625 video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- 626 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
627 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv
628 preprint arXiv:2304.02643*, 2023.
629
- 630 Ilya Kostrikov, Denis Yarats, and Rob Fergus. Image augmentation is all you need: Regularizing
631 deep reinforcement learning from pixels. In *International Conference on Learning Representa-
632 tions*, 2021.
- 633 Nathan Lambert, Brandon Amos, Omry Yadan, and Roberto Calandra. Objective mismatch in
634 model-based reinforcement learning. In *Conference on Learning for Dynamics and Control*, 2020.
635
- 636 Gaspard Lambrechts, Adrien Bolland, and Damien Ernst. Informed POMDP: Leveraging additional
637 information in model-based RL. *Reinforcement Learning Journal*, 1, 2024.
638
- 639 Michael Laskin, Aravind Srinivas, and Pieter Abbeel. Curl: Contrastive unsupervised representa-
640 tions for reinforcement learning. In *International Conference on Machine Learning*, pp. 5639–
641 5650. PMLR, 2020.
- 642 Jessy Lin, Yuqing Du, Olivia Watkins, Danijar Hafner, Pieter Abbeel, Dan Klein, and Anca Dragan.
643 Learning to model the world with language. In *International Conference on Machine Learning*,
644 2024.
- 645 Shakir Mohamed and Danilo Jimenez Rezende. Variational information maximisation for intrinsi-
646 cally motivated reinforcement learning. *Advances in Neural Information Processing Systems*, 28,
647 2015.

- 648 Suraj Nair, Silvio Savarese, and Chelsea Finn. Goal-aware prediction: Learning to model what
649 matters. In *International Conference on Machine Learning*, pp. 7207–7219. PMLR, 2020.
- 650
- 651 Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A univer-
652 sal visual representation for robot manipulation. In *Conference on Robot Learning*, 2022.
- 653
- 654 Tung D Nguyen, Rui Shu, Tuan Pham, Hung Bui, and Stefano Ermon. Temporal predictive coding
655 for model-based planning in latent space. In *International Conference on Machine Learning*, pp.
656 8130–8139. PMLR, 2021.
- 657 Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified,
658 real-time object detection. In *Computer Vision and Pattern Recognition*, pp. 779–788, 2016.
- 659 Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon
660 Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering atari,
661 go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- 662
- 663 Younggyo Seo, Danijar Hafner, Hao Liu, Fangchen Liu, Stephen James, Kimin Lee, and Pieter
664 Abbeel. Masked world models for visual control. In *Conference on Robot Learning*, pp. 1332–
665 1344. PMLR, 2022.
- 666
- 667 John So, Amber Xie, Sunggoo Jung, Jeffrey Edlund, Rohan Thakker, Ali Agha-mohammadi, Pieter
668 Abbeel, and Stephen James. Sim-to-real via sim-to-seg: End-to-end off-road autonomous driving
669 without real data. In *Conference on Robot Learning*, 2022.
- 670 Austin Stone, Oscar Ramirez, Kurt Konolige, and Rico Jonschkowski. The distracting control
671 suite—a challenging benchmark for reinforcement learning from pixels. *arXiv preprint*
672 *arXiv:2101.02722*, 2021.
- 673 Richard S Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *ACM Sigart*
674 *Bulletin*, 2(4):160–163, 1991.
- 675
- 676 Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Bud-
677 den, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. *arXiv*
678 *preprint arXiv:1801.00690*, 2018.
- 679 Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control.
680 In *International Conference on Intelligent Robots and Systems*, pp. 5026–5033. IEEE, 2012.
- 681
- 682 Tongzhou Wang, Simon S Du, Antonio Torralba, Phillip Isola, Amy Zhang, and Yuandong Tian.
683 Denoised mdp: Learning world models better than the world itself. In *International Conference*
684 *on Machine Learning*, 2022.
- 685 Ziyu Wang, Yanjie Ze, Yifei Sun, Zhecheng Yuan, and Huazhe Xu. Generalizable visual reinforce-
686 ment learning with segment anything model. *arXiv preprint arXiv:2312.17116*, 2023.
- 687
- 688 Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Seg-
689 former: Simple and efficient design for semantic segmentation with transformers. In *Advances in*
690 *Neural Information Processing Systems*, 2021.
- 691
- 692 Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Mastering visual continuous control:
693 Improved data-augmented reinforcement learning. *arXiv preprint arXiv:2107.09645*, 2021.
- 694
- 695 Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey
696 Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning.
697 In *Conference on Robot Learning*, 2019.
- 698
- 699 Yuhui Yuan, Lang Huang, Jianyuan Guo, Chao Zhang, Xilin Chen, and Jingdong Wang. Ocnet:
700 Object context network for scene parsing. *arXiv preprint arXiv:1809.00916*, 2018.
- 701
- 702 Amy Zhang, Rowan McAllister, Roberto Calandra, Yarin Gal, and Sergey Levine. Learning invari-
703 ant representations for reinforcement learning without reconstruction. In *International Confer-
704 ence on Learning Representations*, 2021.

702 Marvin Zhang, Sharad Vikram, Laura Smith, Pieter Abbeel, Matthew Johnson, and Sergey Levine.
703 Solar: Deep structured representations for model-based reinforcement learning. In *International*
704 *Conference on Machine Learning*, pp. 7444–7453. PMLR, 2019.

705 Renrui Zhang, Zhengkai Jiang, Ziyu Guo, Shilin Yan, Junting Pan, Hao Dong, Peng Gao, and Hong-
706 sheng Li. Personalize segment anything model with one shot. *arXiv preprint arXiv:2305.03048*,
707 2023.

708 Fangwei Zhong, Kui Wu, Hai Ci, Churan Wang, and Hao Chen. Empowering embodied visual
709 tracking with visual foundation models and offline rl. In *European Conference on Computer*
710 *Vision*, 2024.

711 Chuning Zhu, Max Simchowitz, Siri Gadipudi, and Abhishek Gupta. Repo: Resilient model-based
712 reinforcement learning by regularizing posterior predictability. In *Advances in Neural Information*
713 *Processing Systems*, 2023.

714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A CODE RELEASE

We plan to make the code for Segmentation Dreamer publicly available upon acceptance.

B VISUALIZATION OF TASKS

B.1 DEEPMIND CONTROL SUITE (DMC)

Fig. 6 visualizes the six tasks in DMC (Tassa et al., 2018) used in our experiments. Each row presents the observation from the standard environment, the corresponding observation with added distractions, the ground-truth segmentation mask, and the RGB target with the ground-truth mask applied. Cartpole Swingup Sparse and Cartpole Swingup share the same embodiment and dynamics. Cartpole Swingup Sparse only provides a reward when the pole is upright, whereas Cartpole Swingup continuously provides dense rewards weighted by the proximity of the pole to the upright position. Reacher Easy entails two objects marked with different colors in the segmentation mask, as shown in Fig. 6e 3rd column. Before passing the mask to SD, the mask is converted to a binary format where both objects are marked as *true* as task-relevant.

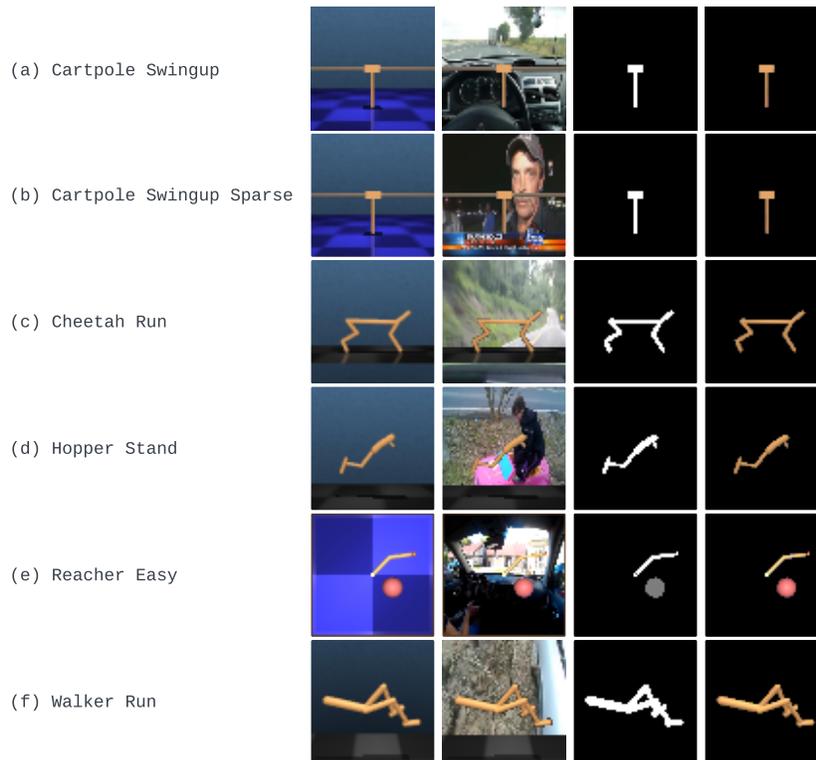


Figure 6: DMC tasks. Left to right: (1) standard environment observations, (2) distracting environment observations, (3) ground-truth segmentation masks, and (4) RGB observations with ground-truth masks applied. We use (4) as auxiliary reconstruction targets in SD^{GT} .

B.2 META-WORLD

Fig. 7 shows the six tasks from Meta-World-V2 used in our experiments. Meta-World is a realistic robotic manipulation benchmark with challenges such as multi-object interactions, small objects, and occlusions.

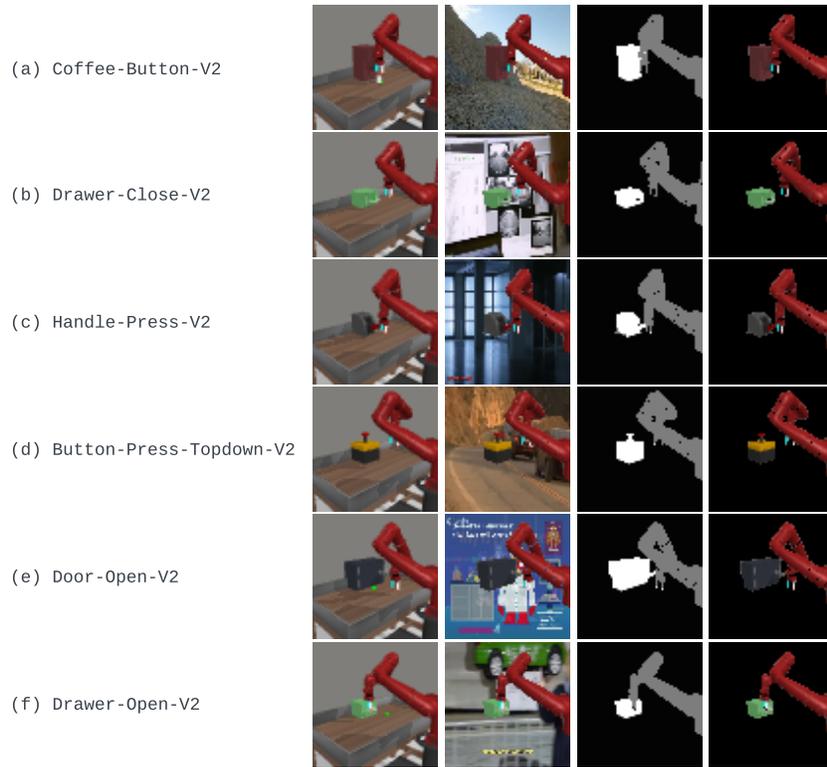


Figure 7: Meta-World tasks. Left to right: (1) standard environment observations, (2) distracting environment observations, (3) ground-truth segmentation masks, and (4) RGB observations with ground-truth masks applied. We use (4) as auxiliary reconstruction targets in SD^{GT} . Masks with multiple classes for different objects are converted to binary masks (all non-background regions are *true* and task-relevant) before use with SD .

C THE IMPACT OF PRIOR KNOWLEDGE

We investigate the impact of accurate prior knowledge of task-relevant objects. Specifically, we conduct additional experiments on Cheetah Run—the task showing the largest disparity between DREAMER* and SD^{GT} in Fig. 3a. In our primary experiment, we designated only the cheetah’s body as the task-relevant object. However, since the cheetah’s dynamics are influenced by ground contact, the ground plate should have also been considered task-relevant.

Fig. 8 (a–c) illustrates the observation with distractions, the auxiliary target without the ground plate, and with the ground plate included, respectively. Fig. 8d compares SD^{GT} trained with different selections of task-relevant objects included in the masked RGB reconstruction targets. We show that including the ground plate leads to faster learning and performance closer to that of the oracle. This highlights the significant influence of prior knowledge on downstream tasks, suggesting that comprehensively including task-relevant objects yields greater benefits.

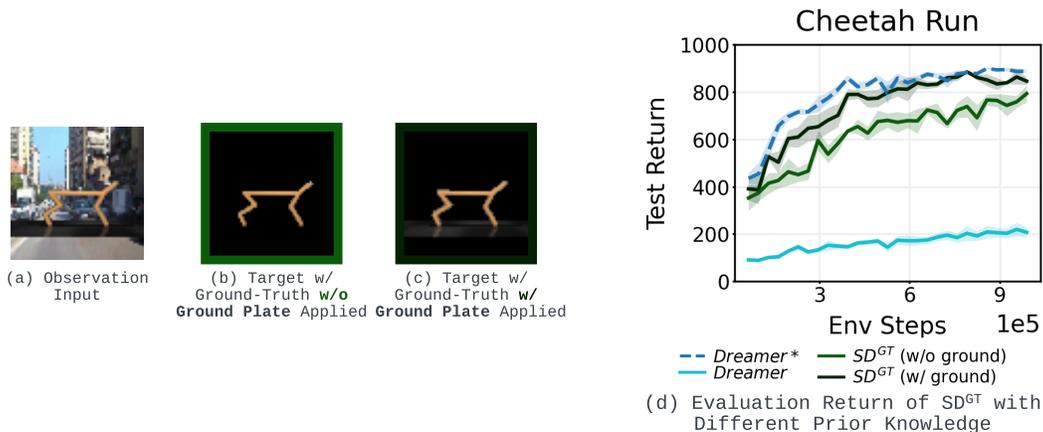


Figure 8: **The impact of prior knowledge on Cheetah Run.** (d) The mean over 4 seeds with the standard error of the mean (SEM) is shaded.

D THE IMPACT OF TEST-TIME SEGMENTATION QUALITY ON PERFORMANCE

We investigate how test-time segmentation quality affects SD^{approx} as well as the *As Input* variation that applies mask predictions to RGB inputs in addition to reconstruction targets. For this analysis, we use PerSAM fine-tuned with a single data point for segmentation prediction. To measure segmentation quality, we compute episodic segmentation quality by averaging over frame-level IoU. In Fig. 9 we plot episode segmentation quality versus test-time reward on the evaluation episodes during the last 10% of training time.

Fig. 9 illustrates that SD^{approx} exhibits greater robustness to test-time segmentation quality compared to the *As Input* variation, with the discrepancy increasing as the IoU decreases. This disparity primarily arises because *As Input* relies on observations restricted by segmentation predictions, and thus its performance deteriorates quickly as the segmentation quality decreases. In contrast, SD^{approx} takes the original observation as input and all feature extraction is handled by the observation encoder, informed by our masked RGB reconstruction objective. Consequently, SD^{approx} maintains resilience to test-time segmentation quality.

An intriguing observation is that a poorly trained agent can lead to poor test-time segmentation quality. For instance, Cartpole Swingup (Sparse) exhibits different segmentation quality distributions between SD^{approx} and *As Input*. This discrepancy occurs because the sub-optimal agent often positions the pole at the cart track edge, causing occlusion and hindering accurate segmentation prediction by PerSAM.

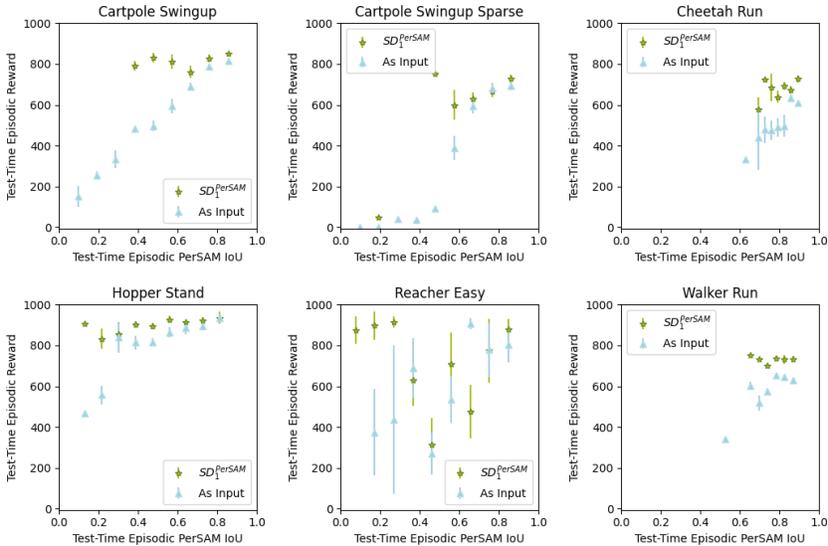


Figure 9: Test-time episodic reward vs PerSAM episodic IoU for SD_1^{PerSAM} and *As Input* (SD_1^{PerSAM} with masked RGB observations as input). SD_1^{PerSAM} is more robust to test-time segmentation prediction errors.

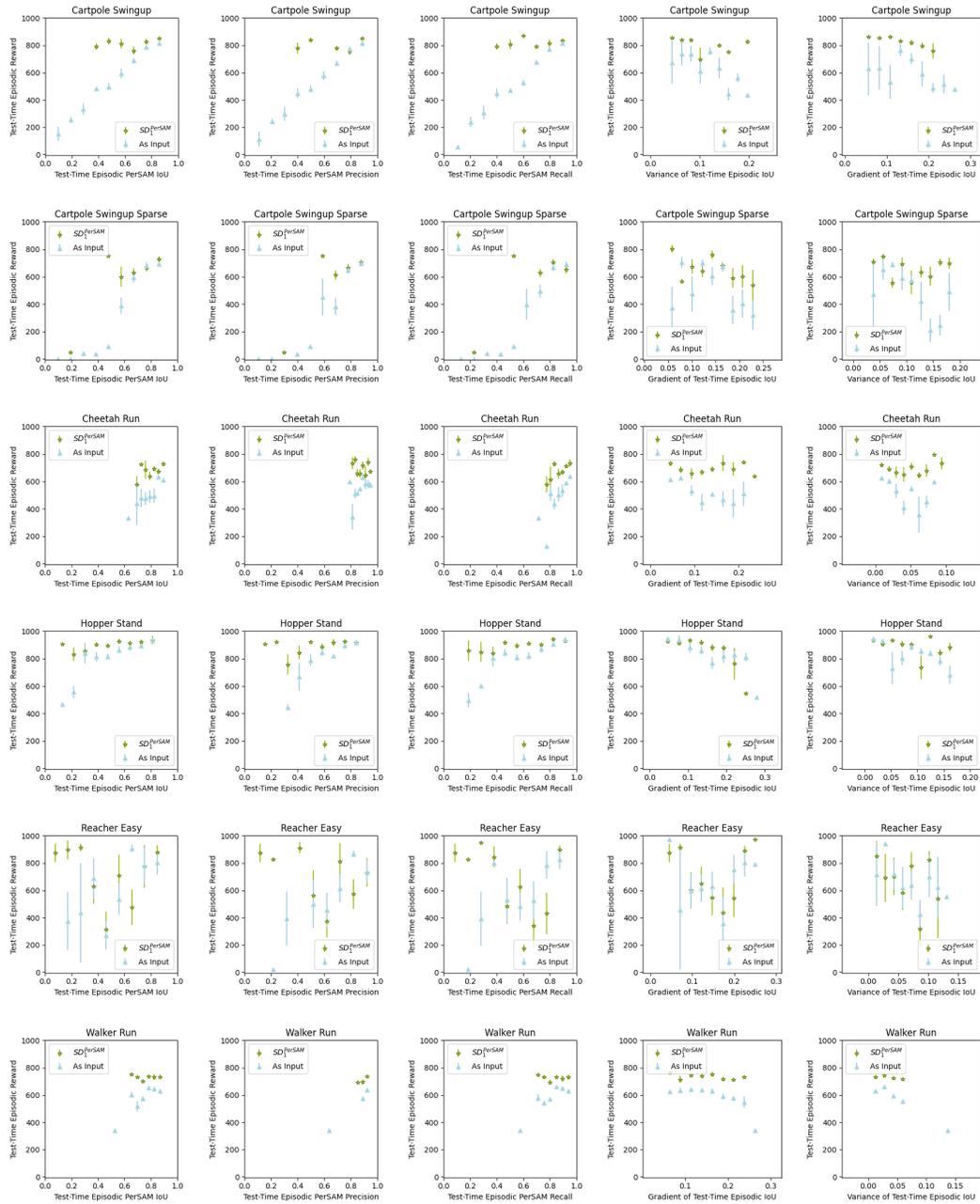


Figure 10: Test-time episodic reward plotted against IoU, precision, recall, IoU variance, and IoU gradient, respectively.

E EXPERIMENTS WITH DIFFERENT TYPES OF DISTRACTIONS

In this section, we investigate how our method performs when faced with types of distractions beyond background distractions. Specifically, we consider three additional types of distractions: foreground distractions, color changes in foreground objects, and camera angle perturbations. These experiments are conducted on the DeepMind Control Suite.

E.1 EXPERIMENT SETUP

To ensure robustness against distractions during testing, we introduce domain randomizations during training. Specifically, both the segmentation models and SD are trained under domain-randomized environments. And, we evaluate our method on a distribution of perturbations that matches the variability introduced at training time.

Foreground Distractions. To simulate distractions that occlude or block task-relevant parts of the scene, we introduce a moving foreground distractor. This is implemented as a blue rectangle rendered near the center of the scene for 4–6 frames every 18–22 frames. These intervals are uniformly sampled each time the distractor appears, meaning approximately 25% of the frames in an episode include the distractor. The distractor moves along pixel-space trajectories defined by randomized Δx and Δy values within the range of $(-3,3)$, which are sampled each time the distractor appears. The goal of introducing this type of distractor is to assess whether our method remains robust in the presence of occlusions that can interfere with task-relevant visual information.

Although this preliminary setup simplifies the distractor’s appearance and trajectory, it can easily be extended to incorporate more complex objects or movement patterns. Given the capabilities of visual foundation models (VFMs), we hypothesize that our method will generalize well to a variety of foreground distractors with different properties.

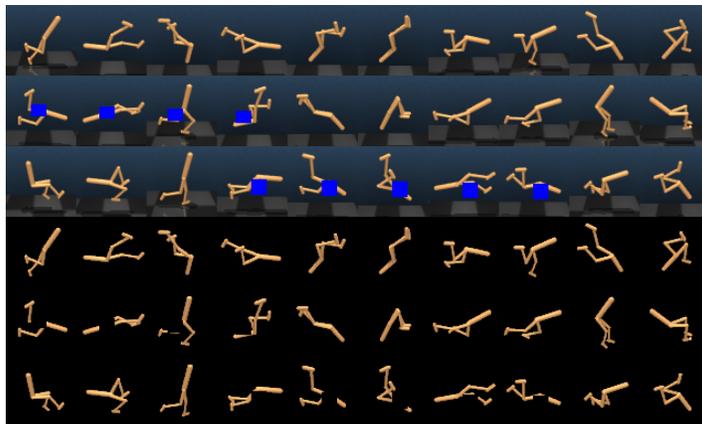


Figure 11: Examples of **foreground distractors** in the environment and corresponding **predictions** from the **segmentation** model that remains **robust** to occlusions in the test set.

Color Changes in Foreground Objects. For color perturbations, we simulate changes in the appearance of the agent or task-relevant objects. Following the approach of Stone et al. (2021), we apply a max delta of 0.1 and set step std to 0.0, resulting in a static color throughout the episode. These changes simulate environmental factors such as lighting variations that may occur during deployment. This experiment evaluates the ability of the model to adapt to changes in the visual characteristics of task-critical elements.

Camera Angle Perturbations. To introduce changes in camera perspective, we follow the implementation of Stone et al. (2021) and apply a scaling factor of 0.1. This results in shifts in the camera view, which simulate real-world deployment scenarios where the agent’s viewpoint may vary due to physical movement or environmental adjustments. These perturbations test the model’s capacity to maintain performance under altered visual perspectives, as illustrated in Fig. 13.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092



Figure 12: Examples of **color** perturbations applied to the agent and corresponding predictions from the segmentation model that remains robust to color changes in the test set.

1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109

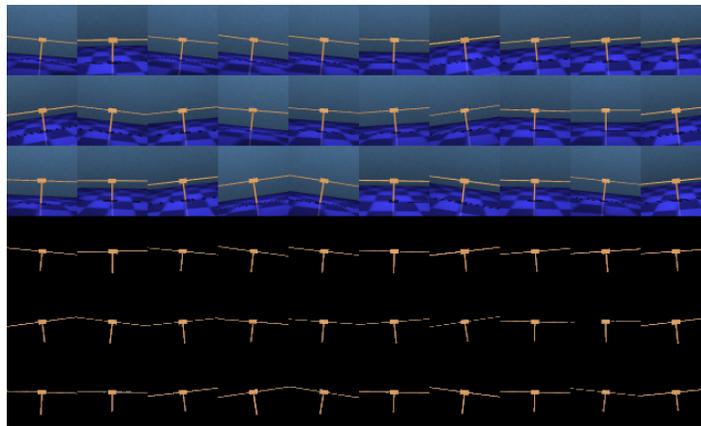


Figure 13: Examples of **camera view** perturbations and corresponding predictions from the segmentation model that remains robust to camera view variations in the test set.

Background Distractions. See Fig. 14 for examples of background perturbations.

1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128

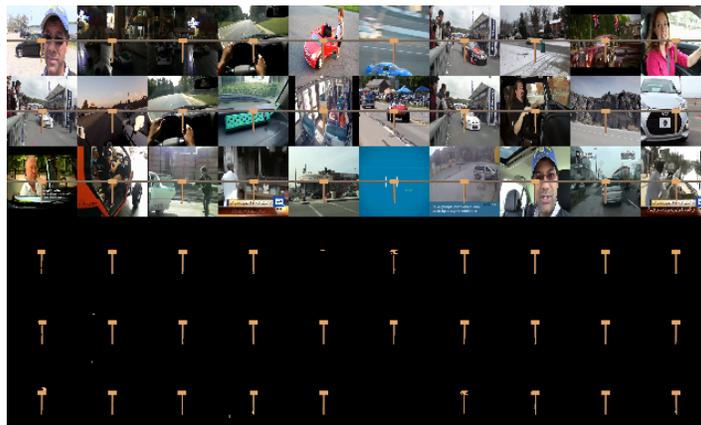
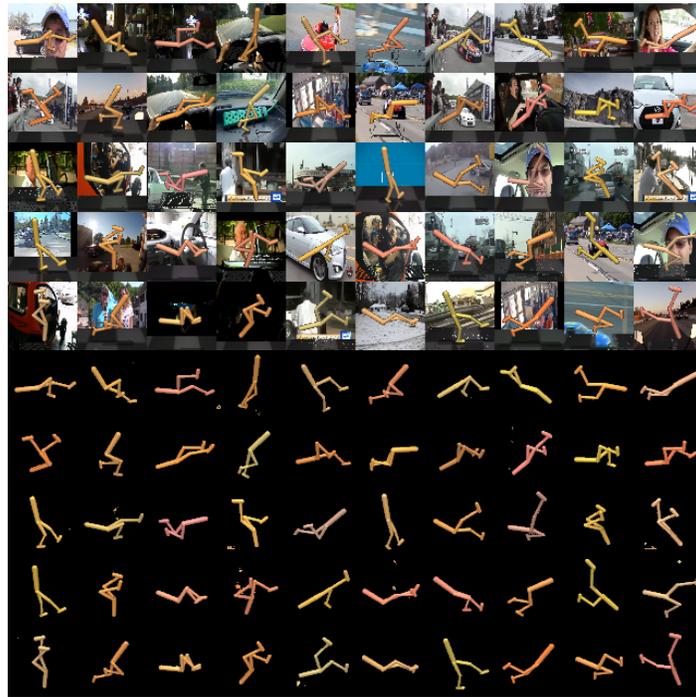


Figure 14: Examples of **background distractions** and predictions from the segmentation model in the test set.

1131
1132
1133

1134 **Background and Color Perturbation.** See Fig. 15 for examples of background and color pertur-
1135 bations.
1136



1159
1160 Figure 15: Examples of **background distractions** and **color** perturbations and predictions from the
1161 segmentation model in the test set.
1162

1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

F SEGMENTATION QUALITY IN META-WORLD

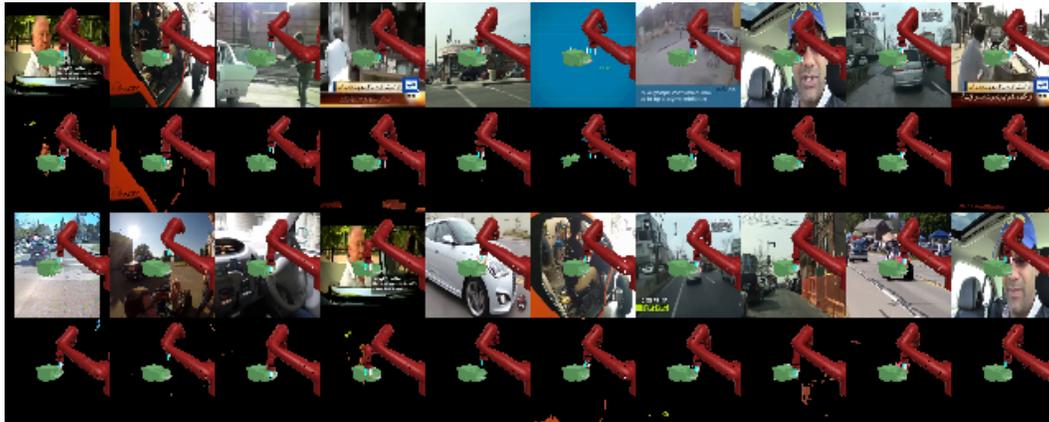


Figure 16: Examples of background perturbations and corresponding **predictions** from the **segmentation** model on Drawer-Open-V2 in the test set.

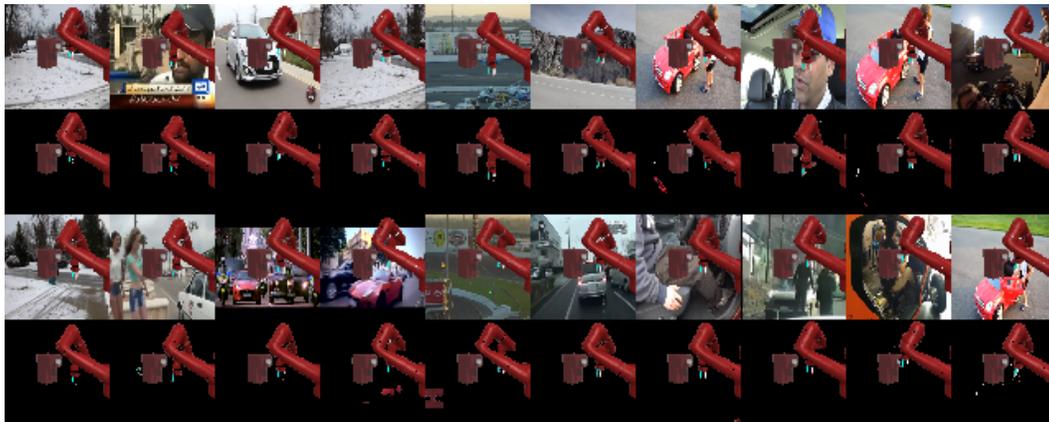


Figure 17: Examples of background perturbations and corresponding **predictions** from the **segmentation** model on Coffee-Button-V2 in the test set.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

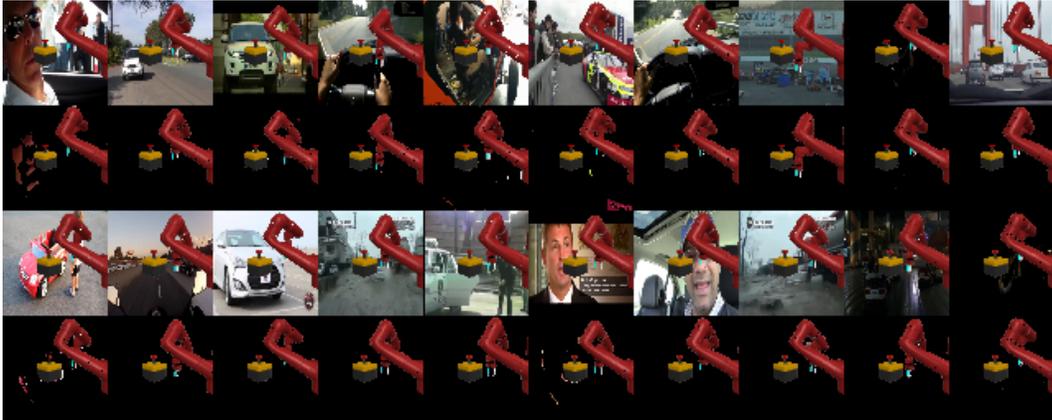


Figure 18: Examples of background perturbations and corresponding **predictions** from the **segmentation** model on Button-Press-Topdown-V2 in the test set.

G ABLATION WITHOUT STOP GRADIENT

Should the SD^{approx.} world model be shielded from gradients of the binary mask decoder head?

To estimate potential regions on RGB targets where task-relevant regions are incorrectly masked out, we train a binary mask prediction head on the world model to help detect false negatives in masks provided by the foundation model. We see better performance when gradients from this binary mask decoder objective are not propagated to the rest of the world model. Thus, the default SD^{approx.} architecture is trained with the gradients of the binary mask branch stopped at its $[h_t; z_t]$ inputs, and the latent representations in the world model are trained only by the task-relevant RGB branch in addition to the standard DREAMER reward/continue prediction and KL-divergence between the dynamics prior and observation encoder posterior. Tab. 2 shows that the performance drops significantly when training without stopping these gradients.

We also examine masks predicted by the binary mask decoder head in Fig. 19. Predictions are coarser grained than their RGB counterparts, lacking details important for predicting intricate forward dynamics. Overall, reconstructing RGB observations with task-relevance masks applied demonstrates itself as a superior inductive bias to learn useful features for downstream tasks compared to binary masks or raw unfiltered RGB observations.

Table 2: Final performance of SD and SD without stop gradient.

Task	SD ₁ ^{PerSAM}	No SG
Cartpole Swingup	730 ± 75	439 ± 81
Cartpole Swingup Sparse	521 ± 92	112 ± 40
Cheetah Run	619 ± 35	376 ± 50
Hopper Stand	846 ± 27	587 ± 127
Reacher Easy	597 ± 97	273 ± 74
Walker Run	730 ± 13	407 ± 62

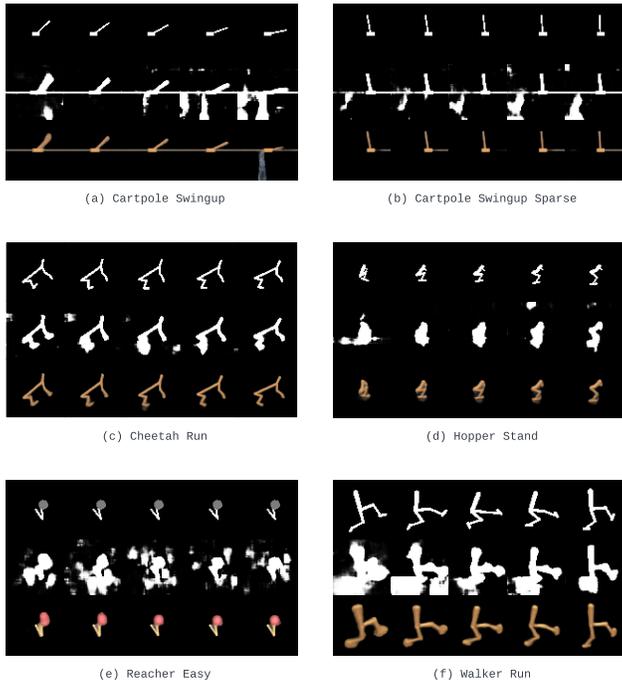


Figure 19: From the top row to the bottom row: (1) ground-truth segmentation masks, (2) SD^{approx.} binary mask predictions, and (3) SD^{approx.} RGB predictions.

H DISTRACTING DMC SETUP

We follow the DBC (Zhang et al., 2021) implementation to replace the background with color videos. The ground plate is also presented in the distracting environment. We used hold-out videos as background for testing. We sampled 100 videos for training from the Kinetics 400 training set of the 'driving car' class, and test-time videos were sampled from the validation set of the same class.

I DISTRACTING META-WORLD SETUP

We test on six tasks from Meta-World-V2. For all tasks, we use the `corner3` camera viewpoint. The maximum episode length for Meta-World tasks is 500 environment steps, with the action repeat of 2 (making 250 policy decision steps). We classify these tasks into *easy*, *medium*, and *difficult* categories based on the training curve of DREAMER* (DREAMER trained in the standard environments). Coffee Button, Drawer Close, and Handle Press are classified as *easy*, and we train baselines on these for 30K environment steps. Button Press Topdown (*medium*) is trained for 100K steps, and Door Open and Drawer Open (*difficult*) are trained for 1M environment steps.

J RESULTS ON META-WORLD WITH SPARSE REWARDS

We also evaluate on sparse reward variations of the distracting Meta-World environments where a reward of 1 is only provided on timesteps when a *success* signal is given by the environment (e.g. objects are at their goal configuration). Rewards are 0 in all other timesteps. The maximum attainable episode reward is 250.

The sparse reward setting is more challenging because the less informative reward signal makes credit assignment more difficult for the RL agent. Fig. 20 shows that our method consistently achieves higher sample efficiency and better performance, showing promise for training agents robust to visual distractions without extensive reward engineering. In Meta-World experiments, TIA (Fu et al., 2021) is not included as it requires exhaustive hyperparameter tuning for new domains and is the lowest-performing method in DMC in general.

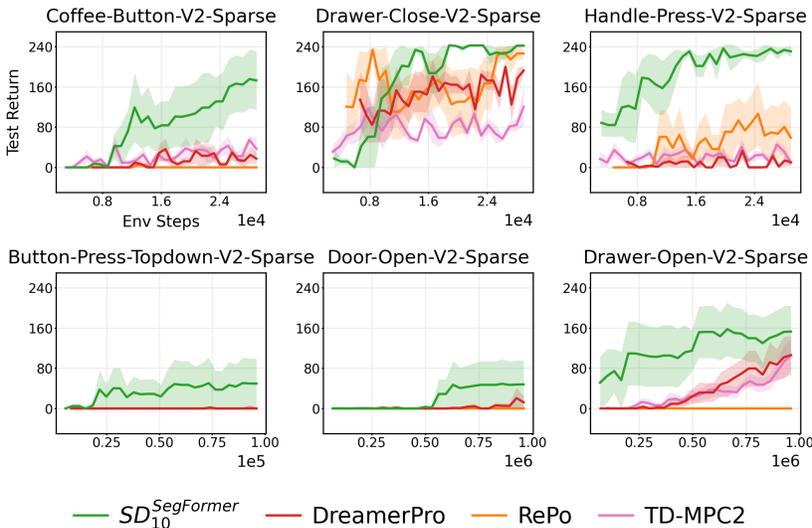


Figure 20: Learning curves on six visual robotic manipulation tasks from Meta-World with sparse rewards.

1404 K FINE-TUNING PERSAM AND SEGFORMER

1405
1406 In this section, we describe how we fine-tune segmentation models and collect RGB and segmenta-
1407 tion mask examples to adapt them.

1408 **PerSAM.** Personalized SAM (PerSAM) (Zhang et al., 2023) is a segmentation model designed
1409 for personalized object segmentation building upon the Segment Anything Model (SAM) (Kirillov
1410 et al., 2023). This model is particularly a good fit for our SD use case since it can obtain a person-
1411 alized segmentation model without additional training by one-shot adapting to a *single* in-domain
1412 image. In our experiments, we use the model with ViT-T as a backbone.

1413 **SegFormer.** We use 5 or 10 pairs of examples to fine-tune SegFormer (Xie et al., 2021) MiT-b0.

1414 To collect a one-shot in-domain RGB image and mask example for DMC and MetaWorld exper-
1415 iments, we sample a state from the initial distribution p_0 and render the RGB observation. In a
1416 few-shot scenario, we deploy a random agent in to collect more diverse observations from more
1417 diverse states.

1418 To generate the associated masks for these states, we make additional queries to the simulation
1419 rendering API. We represent the pixel values for background and irrelevant objects as *false* and
1420 task-relevant objects as *true*. In multi-object cases, we may perform a separate adaptation operation
1421 for each task-relevant object, resulting in more than 2 mask classes. In such cases, before integrating
1422 masks with SD^{approx} , we will combine the union of the mask classes for all pertinent objects as a
1423 single *true* task-relevant class, creating a binary segmentation mask compatible with our method.

1424 In cases where example masks cannot be programmatically extracted, because such a small number
1425 of examples are required (1-10), it should also be very feasible for a human to use software to
1426 manually annotate the needed mask examples from collected RGB images.

1429 L DETAILS ON SELECTIVE L_2 LOSS

1430 The binary mask prediction branch in SD^{approx} is equipped with the sigmoid layer at its output. In
1431 order to obtain binary mask_{SD}, we binarize the SD binary mask prediction with a threshold of 0.9.

1435 M DETAILS ON BASELINES

1436 It is known that RePo (Zhu et al., 2023) outperforms many earlier works (Fu et al., 2021; Hansen
1437 et al., 2022; Zhang et al., 2021; Wang et al., 2022; Gelada et al., 2019) and that DreamerPro (Deng
1438 et al., 2022) surpasses TPC (Nguyen et al., 2021). However, these two groups of works have been
1439 using slightly different environment setups and have not been compared with each other despite
1440 addressing the same high-level problem on the same DMC environments. In our experiments, we
1441 evaluate the representatives in each cluster on a common ground (See Appendix H) and compare
1442 them with our method.

1443 In our experiments, we use hyperparameters used in the original papers for all the baselines, ex-
1444 cept RePo (Zhu et al., 2023) in Meta-World. RePo does not have experiments on Meta-World in
1445 which case we use hyperparameters used for Maniskill2 (Gu et al., 2023) which is another robot
1446 manipulation benchmark.

1449 N EXTENDED RELATED WORK

1450 There are several model-based RL approaches which introduce new auxiliary tasks. Dynalang (Lin
1451 et al., 2024) integrates language modeling as a self-supervised learning objective in world-model
1452 training. It shows impressive performance on benchmarks where the dynamics can be effectively
1453 described in natural language. However, it is not trivial to apply this method in low-level control sce-
1454 narios such as locomotion control in DMC. Informed POMDP (Lambrechts et al., 2024) introduces
1455 an information decoder which uses privileged simulator information to decode a sufficient statistic
1456 for optimal control. This shares an idea of using additional information available at training time
1457 with our method SD^{GT} . Although this can be effective on training in simulation where well-shaped

1458 proprioceptive states exist, it cannot be applied to cases where such information is hard to obtain. In
1459 goal-conditioned RL, GAP (Nair et al., 2020) proposed to decode the difference between the future
1460 state and the goal state to help learn goal-relevant features in the state space.

1462 O LIMITATIONS

1464 Segmentation Dreamer achieves excellent performance across diverse tasks in the presence of dis-
1465 tractions and provides a human interface to indicate task relevance. This capability enables prac-
1466 titioners to readily train an agent for their specific purposes without suffering from poor learning
1467 performance due to visual distractions. However, there are several limitations to consider.

1468 First, since SD^{approx} harnesses a segmentation model, it can become confused when a scene contains
1469 distractor objects that resemble task-relevant objects. This challenge can be mitigated by combining
1470 our method with approaches such as InfoPower (Bharadhwaj et al., 2022), which learns control-
1471 lable representations through empowerment (Mohamed & Jimenez Rezende, 2015). This integra-
1472 tion would help distinguish controllable task-relevant objects from those with similar appearances
1473 but move without agent interaction.

1474 Second, our method does not explicitly address randomization in the visual appearance of *task-*
1475 *relevant* objects, such as variations in brightness, illumination, or color. Two observations of the
1476 same internal state but with differently colored task-relevant objects may be guided toward differ-
1477 ent latent representations because our task-relevant "pixel-value" reconstruction loss forces them to
1478 be differentiated. Ideally, these observations should map to the same state abstraction since they
1479 exhibit similar behaviors in terms of the downstream task. Given that training with pixel-value
1480 perturbations on task-relevant objects is easier compared to dealing with dominating background
1481 distractors (Stone et al., 2021), our method is expected to manage such perturbations effectively
1482 without modifications. However, augmenting our approach with additional auxiliary tasks based on
1483 behavior similarity (Zhang et al., 2021) would further enhance representation learning and directly
1484 address this issue.

1485 Finally, our approximation model faces scalability challenges when task-relevant objects constitute
1486 an open set. For instance, in autonomous driving scenarios, obstacles are task-relevant but cannot
1487 be explicitly specified. While our method serves as an effective solution when task-relevant objects
1488 are easily identifiable, complementary approaches should be considered when this assumption does
1489 not hold true.

1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511