LEARNING OBJECT-CENTRIC LATENT DYNAMICS FOR REINFORCEMENT LEARNING FROM PIXELS

Anonymous authors

Paper under double-blind review

Abstract

Learning a latent dynamics model provides a task-agnostic representation of an agent's understanding of its environment. Leveraging this knowledge for modelbased reinforcement learning holds the potential to improve sample efficiency over model-free methods by learning inside imagined rollouts. Furthermore, because the latent space serves as input to behavior models, the informative representations learned by the world model facilitate efficient learning of desired skills. Most existing methods rely on holistic representations of the environment's state. In contrast, humans reason about objects and their interactions, forecasting how actions will affect specific parts of their surroundings. Inspired by this, we propose Slot-Attention for Object-centric Latent Dynamics (SOLD), a novel algorithm that learns object-centric dynamics models in an unsupervised manner from pixel inputs. We demonstrate that the structured latent space not only improves model interpretability but also provides a valuable input space for behavior models to reason over. Our results show that SOLD outperforms DreamerV3, a state-of-theart model-based RL algorithm, across a range of benchmark robotic environments that evaluate for both relational reasoning and low-level manipulation capabilities.

026 027 028

029

000

001

002 003 004

005

010 011

012

013

014

015

016

017

018

019

020

021

022

023

1 INTRODUCTION

Advances in reinforcement learning (RL) have showcased the ability to learn sophisticated control strategies through interaction, achieving superhuman performance in domains ranging from board games (Silver et al., 2016) to drone racing (Kaufmann et al., 2023). While these approaches excel in settings where explicit models of the environment are available or abundant data can be collected, learning complex control tasks in a sample-efficient manner remains a significant challenge. Modelbased RL (MBRL) has emerged as a promising approach to address this limitation by constructing models of the environment's dynamics. Notably, the Dreamer framework (Hafner et al., 2019; 2020; 2023) has demonstrated improved sample efficiency over model-free baselines by learning behaviors solely through imagined rollouts.

While these research efforts have produced world models capable of accurately predicting the dy-039 namics of visual tasks, they rely on a holistic representation of the environment's state. In contrast, 040 humans perceive the world by parsing scenes into individual objects (Spelke, 1990), anticipating 041 how their actions will influence specific components of their surroundings. This ability is crucial in 042 complex tasks that require reasoning about multiple objects and their interactions, which is common 043 in robotic manipulation. Learning structured representations of an environment through objects introduces a powerful inductive bias. The learned representations not only enhance the interpretability 045 of the dynamics prediction but also foster decision-making by providing a structured input space for behavior models to reason over. Despite these advantages, the integration of object-centric repre-046 sentations and world models remains largely underexplored. To the best of our knowledge, no prior 047 work has introduced a method that performs object-centric model-based RL directly from pixels. 048

To address the limitations of holistic representations in model-based RL, we propose Slot-Attention for Object-centric Latent Dynamics (SOLD), a novel algorithm that leverages structured, object-centric states within the latent space of its world model. Our method introduces two key innovations to the model-based RL framework. The first contribution is a dynamics model that predicts future frames in terms of their slot representation. To achieve this, we extend OCVP (Villar-Corrales et al., 2023) into an action-conditional model by projecting action commands into the slot



Figure 1: We evaluate SOLD on diverse visual environments from (left to right) our multi-object robotic control benchmark, Meta-World (Yu et al., 2019), and DM-Control (Tassa et al., 2018).

071 dimension. A transformer-based (Vaswani et al., 2017) transition model processes the set of object-072 slots and actions to predict the subsequent frame. Notably, the dynamics model is trained solely 073 from pixels through a loss on the reconstructions and slot representations of the predicted frames. 074 The second key component is a transformer-based model backbone, which we call *Slot Aggrega*tion Transformer, that is used to make time-step-wise predictions from the history of object slots. 075 Specifically, it serves as the architectural backbone to the reward, value, and action model, allowing 076 us to perform model-based RL training on the basis of object-centric representations. 077

078 For systematic evaluation, we introduce a suite of visual robotics tasks, shown in Figure 1, that require varying levels of relational reasoning and manipulation capabilities. We perform an extensive 080 comparison on this benchmark, demonstrating that our method achieves superior performance to the state-of-the-art MBRL algorithm DreamerV3 (Hafner et al., 2023). Furthermore, we apply SOLD 081 to tasks from two RL benchmarks that were not designed to be object-centric, providing evidence of 082 the generalizability of our framework. In summary, we make the following contributions: 083

- We introduce SOLD, which is to the best of our knowledge, the first object-centric modelbased RL algorithm that learns entirely from pixel inputs.
- Our method outperforms DreamerV3 across a range of visual robotics environments, excelling in tasks that require both relational reasoning and low-level manipulation skills.
- We overcome limitations of prior object-centric methods in RL. Namely, we show that our object-centric encoder-decoder module can be adapted to state distributions vastly different from those seen under a random policy during pre-training – an essential capability for solving many complex tasks. Additionally, we show its generalization potential to environments not explicitly designed for object-centric reasoning.

2 BACKGROUND

097 Slot Attention for Video (SAVi) SOLD employs SAVi (Kipf et al., 2022), an encoder-decoder architecture with a structured bottleneck composed of N permutation-equivariant object embeddings denoted as slots, in order to recursively parse a sequence of video frames $o_{0:\tau} = \{o_0, ..., o_{\tau}\}$ into their object representations $Z_{0:\tau} = \{Z_0, ..., Z_{\tau}\}, Z_t \in \mathbb{R}^{N \times D_Z}$. At time t, SAVi encodes the input video frame o_t into a set of feature maps $F_t \in \mathbb{R}^{L \times D_h}$, where L is the size of the flattened grid 100 101 (i.e. $L = \text{width} \cdot \text{height}$), and uses Slot Attention (Locatello et al., 2020) to iteratively refine the 102 previous slot representations conditioned on the current features. Slot Attention performs cross-103 attention between the slots and image features with the attention coefficients normalized over the 104 slot dimension, thus encouraging the slots to compete to represent feature locations: 105

106 107

084

085

087

089

092

093 094 095

$$\boldsymbol{A} \doteq \operatorname{softmax}_{N}\left(\frac{q(\boldsymbol{Z}_{t-1}) \cdot k(\boldsymbol{F}_{t})^{T}}{\sqrt{D}}\right) \in \mathbb{R}^{N \times L},\tag{1}$$



121(a) World Model Learning: SAVi encodes images o_t 122into slots Z_t , which are predicted by the dynamics123model given history of slots and actions. We recon-124resentation to shape the dynamics prediction.

(b) Behavior Learning: The actor and critic are trained via imagined rollouts in the latent space of the world model. Trajectories start after S seed frames (visualized for S = 1) and predict forward with actions sampled from the actor network.

Figure 2: *SOLD* is trained by concurrently making the world model consistent with replayed experiences and learning behaviors through latent imagination.

where q and k are learned linear mappings to a common dimension D. The slots are then independently updated via a shared Gated Recurrent Unit (Cho, 2014) (GRU) followed by a residual MLP:

$$\boldsymbol{Z}_{t} \doteq \mathrm{MLP}(\mathrm{GRU}(\boldsymbol{A} \cdot \boldsymbol{v}(\boldsymbol{F}_{t}), \boldsymbol{Z}_{t-1})) \text{ with } \boldsymbol{A}_{n,l} \doteq \frac{\boldsymbol{A}_{n,l}}{\sum_{i=0}^{L-1} \boldsymbol{A}_{n,i}},$$
(2)

where v is a learned linear projection. The steps described in Equations 1 and 2 can be repeated multiple times with shared weights to iteratively refine the slots and obtain an accurate object-centric representation of the scene. Finally, SAVi independently decodes each slot of Z_t into per-object images and alpha masks, which can be normalized and combined via weighted sum to render video frames. SAVi is trained end-to-end in a self-supervised manner with an image reconstruction loss.

140 **OCVP** Our dynamics model builds on OCVP (Villar-Corrales et al., 2023) in order to autore-141 gressively predict future object slots conditioned on past object states. OCVP is a transformer-142 encoder model that decouples the processing of object dynamics and interactions, thus leading to 143 interpretable and temporally consistent object predictions while retaining the inherent permutationequivariant property of the object slots. This is achieved through the use of two specialized self-144 attention variants: temporal attention updates a slot representation by aggregating information from 145 the corresponding slot up to the current time step, without modeling interactions between distinct 146 objects, whereas relational attention models object interactions by jointly processing all slots from 147 the same time step. 148

149

125

126

127 128 129

130

150 151

3 SLOT-ATTENTION FOR OBJECT-CENTRIC LATENT DYNAMICS

We propose *Slot-Attention for Object-centric Latent Dynamics (SOLD)*, a method that combines model-based RL akin to the Dreamer framework (Hafner et al., 2019; 2020; 2023) with objectcentric representations. The three core components of our algorithm are: the *object-centric world model*, which predicts the effects of actions on the environment, the *critic*, which estimates the value of a given state, and the *actor* that selects actions to maximize this value.

157 Figure 2 gives an overview of the training process. The world model operates on structured latent 158 states by splitting the environment into its constituent objects and then composing future frames 159 via the predicted states of these individual components. Specifically, we pretrain a SAVi encoder-160 decoder model (Kipf et al., 2022) on random sequences from the environment to extract object-161 centric representations. After pretraining, all components of the world model are trained jointly 162 using replayed experiences from the agent's interaction with the environment. The actor and critic are trained on imagined sequences of structured latent states. We execute actions sampled from the
 actor model in the environment and append the resulting experiences to the replay buffer. Detailed
 explanations of world model learning and behavior learning are provided in Sections 3.1 and 3.2,
 respectively.

167 3.1 WORLD MODEL LEARNING

166

World models compress an agent's experience into a predictive model that forecasts the outcomes of potential actions. By simulating rollouts within the internal model, agents can learn desired behaviors in a sample-efficient manner. When the inputs are high-dimensional images, it is helpful to learn compact state representations, enabling prediction within this latent space. This type of model, called latent dynamics model, allows for efficient prediction of many latent sequences in parallel.

Most prior works rely on generating a single, holistic representation of the environment's state, which contrasts with findings from cognitive psychology. Humans perceive scenes as compositions of objects (Spelke, 1990) and reason about how their actions affect distinct parts of their environment. Furthermore, environment dynamics can be compactly explained in terms of objects and their interactions Battaglia et al. (2016). Therefore, we propose to structure the latent space by decomposing visual environments into their constituent parts.

Components To create a world model that operates on object-centric latent representations, we build on top of OCVP (Villar-Corrales et al., 2023). We begin by pretraining SAVi on a dataset of 10⁶ frames from random episodes. Having a sufficiently large initial dataset is crucial for meaningful object-centric representations to emerge. We do not freeze the pretrained encoder-decoder models, allowing slots to adapt to novel configurations that do not occur during random pre-training. The sequence of object slots $Z_{0:t}$ alongside the action commands $a_{0:t}$ serve as inputs to our transformerbased dynamics model which predicts the slot representation of the next frame \hat{Z}_{t+1} :

Encoder:	$oldsymbol{Z}_t = e_\eta(oldsymbol{o}_t),$	
Decoder:	$\hat{\boldsymbol{o}}_t = d_\eta(\boldsymbol{Z}_t),$	(2)
Dynamics model:	$\hat{oldsymbol{Z}}_{t+1} = p_{\psi}(oldsymbol{Z}_{0:t},oldsymbol{a}_{0:t}),$ and	(3)
Reward predictor:	$\hat{r}_t \sim p_{\zeta}(\hat{r}_t \mid \boldsymbol{Z}_{0:t}).$	

Object-centric dynamics learning For the dynamics model, we follow the sequential attention pattern proposed in Villar-Corrales et al. (2023), which disentangles relational and temporal attention to decouple the processing object dynamics and interactions. During training, we provide the slot representation of S seed frames as context. We append the predictions to the context and apply this process in an autoregressive manner to predict the subsequent T frames. We do not employ teacher forcing so that the dynamics model learns to handle its own imperfect predictions. To shape the predicted representations, we reconstruct the subsequent frame \hat{o}_{t+1} and extract the SAVi representations of the actual frame Z_{t+1} to compute the hybrid dynamics loss:

$$\mathcal{L}_{dyn}(\psi) \doteq \sum_{t=S}^{S+T-1} \left[\underbrace{\left\| \hat{\boldsymbol{Z}}_t - \boldsymbol{e}_{\eta}(\boldsymbol{o}_t) \right\|_2^2}_{\text{Joint embedding}} + \underbrace{\left\| \hat{\boldsymbol{o}}_t - \boldsymbol{o}_t \right\|_2^2}_{\text{Reconstruction}} \right].$$
(4)

204 205

207

208

202 203

For all loss terms, we specify the parameter group that is being optimized and omit stop-gradients for other models to avoid cluttering the notation.

Reward model learning The reward predictor solves a regression problem where the prediction depends on the slot representations but is not directly tied to any single slot. To address this, we introduce the *Slot Aggregation Transformer (SAT)* as an architectural backbone, which introduces output tokens and a variable number of register tokens for all time-steps. Register tokens, recently shown to enhance computation in vision transformers (Darcet et al., 2024), can aid computation when processing a set of inputs to produce a singular output. To encode the position information, we adopt ALiBi (Press et al., 2022) in place of absolute position encoding. ALiBi introduces linear biases directly into the attention scores, effectively encoding token recency. This approach helps to

216 generalize to sequences longer than those seen during training. A detailed description of the SAT 217 can be found in Section C.3. To efficiently represent a wide range of reward values, we avoid directly 218 predicting a scalar reward. Instead, the MLP head f_{ζ} outputs logits of a softmax distribution over 219 K exponentially spaced bins b_i . The predicted reward can then be computed as the expectation over 220 these bins:

$$\boldsymbol{b} \doteq \operatorname{symexp}([-20, ..., +20]), \qquad \hat{r}_t \doteq \operatorname{softmax}(f_{\zeta}^{\mathsf{MLP}}(\boldsymbol{h}_t))^T \boldsymbol{b}, \tag{5}$$

where h_t are the output tokens after being processed by the SAT backbone. To formulate the loss, the true reward r_t is first transformed using the symlog function (Webber, 2012) and then encoded via a two-hot encoding strategy (Bellemare et al., 2017; Schrittwieser et al., 2020). The model is trained to maximize the log-likelihood of the two-hot encoded reward distribution under the predicted distribution:

221

238

 $\mathcal{L}_{\text{rew}}(\zeta) \doteq -\sum_{t=0}^{T-1} \log p_{\zeta}(r_t \mid \mathbf{Z}_{0:t}).$ (6)

3.2 BEHAVIOR LEARNING

Our strategy of using the world model for behavior learning builds upon the Dreamer framework. At the core of this method lies the process of latent imagination, visualized in Figure 2b, which trains the actor and critic networks purely on imagined trajectories predicted by the world model. Since both the actor and critic operate on the latent state, they benefit from the structured representation learned by the world model. The architecture of both models mirrors that of the reward predictor, consisting of a SAT backbone that processes the slot histories followed by an MLP head:

Actor:
$$\boldsymbol{a}_t \sim \pi_{\theta}(\boldsymbol{a}_t \mid \boldsymbol{Z}_{0:t}),$$
 Critic: $\hat{R}_t \doteq \mathbb{E}[v_{\phi}(\hat{R}_t \mid \boldsymbol{Z}_{0:t})].$ (7)

Critic learning To account for rewards beyond the imagination horizon T = 15, the critic is trained to estimate the expected return under the current actor's behavior. Since no ground truth is available for these estimates, we compute bootstrapped λ -returns (Sutton & Barto, 2018), R^{λ} , via temporal difference learning. These returns integrate predicted rewards \hat{r} and values \hat{R} to form the target for the value model:

$$R_t^{\lambda} \doteq \hat{r}_{t+1} + \gamma \left((1-\gamma)\hat{R}_{t+1} + \lambda R_{t+1}^{\lambda} \right) \quad \text{where} \quad R_T^{\lambda} \doteq \hat{R}_T, \tag{8}$$

which is trained to minimize the resulting loss:

246 247

245

248

249

250 251

264 265 266 $\mathcal{L}_{\text{critic}}(\phi) \doteq -\sum_{t=0}^{T-1} \log v_{\phi}(R_t^{\lambda} \mid \mathbf{Z}_{0:t}).$ (9)

We decouple the gradient scale from value prediction through same approach as in the reward model, predicting a categorical distribution over exponentially spaced bins. To stabilize learning, we regularize the critic's predictions towards the outputs of an exponentially moving average (EMA) of its own parameters (Mnih et al., 2015; Hafner et al., 2023).

256 Actor learning The actor is optimized to select actions that maximize its expected return while 257 encouraging exploration through an entropy regularizer. Its model architecture is similar to the 258 critic and reward predictor, but instead of regressing a scalar value, it predicts the parameters of 259 the action distribution. Specifically, the MLP head outputs the mean μ_t and standard deviation σ_t to parameterize a normal distribution $\mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\sigma}_t | \boldsymbol{Z}_{0:t})$ over possible actions. The trade-off in the 260 actor's loss function weights expected returns with maintaining randomness in the actor outputs and 261 is hence subject to reward scale and frequency of the current environment. To adapt to varying scales 262 of value estimates across different environments, we use a normalization factor scale_V: 263

$$\mathcal{L}_{actor}(\theta) \doteq -\sum_{t=0}^{T-1} \frac{\hat{R}_t^{\lambda}}{\max(1, \text{scale}_V)} + \eta H(\mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\sigma}_t \mid \boldsymbol{Z}_{0:t})),$$
(10)

where the value normalization is computed via the EMA of the 5th and 95th percentile of the value estimates (Hafner et al., 2023):

$$\operatorname{scale}_{V} \doteq \operatorname{EMA}\left(\operatorname{Per}(\hat{R}_{t}^{\lambda},95) - \operatorname{Per}(\hat{R}_{t}^{\lambda},5),0.99\right).$$
(11)



Figure 3: Open-loop predictions on the tasks *PickAndPlace-Distinct* (top) and *Hammer* (bottom). Starting from a single context frame, our model predicts the next 50 frames by propagating the individual slot representations forward without access to any intermediate images.

RESULTS

We evaluate SOLD on a suite of robotic manipulation tasks, designed to test for both complex relational reasoning as well as low-level manipulation capabilities. Further, we apply our method to environments that are not designed as object-centric tasks. Specifically, we test on two environments from Meta-World (Yu et al., 2019) and DM-Control (Tassa et al., 2018), respectively.

Baselines We design the experiments to achieve two main objectives: first, to specifically assess the impact of the object-centric paradigm in our method by comparing it to a baseline that replaces the object-centric encoder-decoder modules with a standard convolutional architecture (Ours w/o OCE); and second, to evaluate our approach against the best available competitor from the literature by benchmarking it against DreamerV3 (Hafner et al., 2023), a state-of-the-art MBRL algorithm known for its strong performance across a wide range of tasks. Here, we choose the 12 million parameter version, to match the parameter count of our own model. Further details about the baselines are provided in Appendix D.

Environments We introduce a suite of eight object-centric robotic control environments designed to test both relational reasoning and manipulation capabilities. These environments feature two types of problems: *Reach* tasks, where the agent must identify a target and move the end-effector to its location, and manipulation tasks (Push and PickAndPlace), where the agent identifies a target block and moves it to a designated goal. The action-space is 4-dimensional, where the first three components represent the desired movement direction of the end effector, and the fourth controls the gripper. On the *Reach* and *Push* tasks, commands to the gripper are ignored, with the gripper fixed in a closed configuration, as gripping is not required to solve these tasks. To test varying levels of relational reasoning difficulty, we design the following tasks:

- *Specific* The target object is red, with 0 to 4 distractor objects of random, distinct colors present in the scene.



Figure 4: Final success rates across the eight evaluated environments. Overall, SOLD compares favorably to DreamerV3, showing slight improvements on the Specific tasks and a significant performance boost on the Distinct variants.

• Distinct Inspired by the odd-one-out task in cognitive science (Crutch et al., 2009; Beatty & Vartanian, 2015), this task presents 3 to 5 objects, and the target is the one that differs in color from all the others.

For the *Reach* task, we also introduce two more challenging variants:

- Specific-Relative Instead of reaching for the red object, the goal is to reach the reddest object, determined by the perceptual CIEDE2000 (Sharma et al., 2005) distance.
- *Distinct-Groups* The environment always contains 5 targets, and the goal is to reach the one that appears only once.

On these two additional reach tasks, we reuse the SAVi models that were pre-trained for Reach-350 Specific and Reach-Distinct, respectively without modification. Further details about these environments are provided in Appendix E.

352 In addition to these tasks, we evaluate our approach on environments not originally designed for 353 object-centric learning to investigate its generalizability. We include the Button-Press and Hammer 354 tasks from the Meta-World benchmark (Yu et al., 2019), both featuring objects with complex shapes 355 and textures, testing the model's ability to handle more diverse visual inputs. Finally, we assess our 356 method on the Cartpole-Balance and Finger-Spin environments from the DM-Control suite (Tassa 357 et al., 2018), which represent vastly different domains not typically associated with object-centric approaches. An overview of all studied environments is given in Figure 1. 358

359 360

361

333

334

335 336 337

338

339

340 341

342 343

344

345 346

347

348 349

351

4.1 **OBJECT-CENTRIC DYNAMICS LEARNING**

362 The object-centric representations learned by SAVi can be seen in the context-frames of Figure 3. 363 The slots effectively decompose the visual scene, with most slots representing distinct objects and 364 three slots capturing different parts of the respective robots. This part-whole segmentation demon-365 strates that the slots can meaningfully identify separate parts of a larger object, representing the 366 gripper jaws in the first example and different parts of the kinematic chain of the Sawyer robot in 367 the second. Notably, the sharp mask predictions show that each slot isolates information about the 368 specific object it represents (see also Section F). This property is crucial for object-centric behavior 369 learning, as it enables subsequent components to reason about task-relevant objects while ignoring irrelevant information. 370

371 Further, the open-loop predictions visualized in Figure 3 demonstrate that the object decomposition 372 learned by the SAVi model is preserved throughout the prediction sequence. On the *PickAndPlace*-373 Distinct task, the movement of the robot and the blocks is predicted with high accuracy, showcasing 374 the model's ability to capture complex physical interactions. Furthermore, the model effectively 375 handles occlusions, as evidenced by the continued precise prediction of the spherical red target. In the second example, the slots capture the intricate shape of the hammer and nail-box. The predictions 376 remain reliable over a long horizon, even during interactions between the robot and the hammer, and 377 between the hammer and nail.



Figure 5: Achieved return over the training duration on our eight benchmark environments. The dotted vertical line indicates the offset of our method to account for the data used during pre-training.

4.2 BEHAVIOR LEARNING

To assess SOLD's performance across our suite of robotic control tasks, we train each method three 400 times with different random seeds on each environment. The final success rate achieved by each 401 method is shown in Figure 4. SOLD consistently outperforms the non-object-centric baseline, often 402 by a significant margin, highlighting the effectiveness of object-centric representations for these 403 problems. While the non-object-centric baseline is able to make progress on the Specific variants of 404 the tasks, it struggles to perform the relational reasoning required to solve the *Distinct* versions. As 405 a result, we decided excluded it from our evaluation of the two more advanced versions of the Reach 406 tasks. When compared to the state-of-the-art method DreamerV3, our model shows competitive or superior performance. This advantage is particularly pronounced for tasks that require relational 407 reasoning between objects, namely the Distinct variants and the Specific-Relative task. 408

409 These results confirm our hypothesis that learning a structured latent representation in the world 410 model benefits downstream tasks that require object reasoning, which are common in robotics con-411 texts. Beyond that, when examining the performance over the course of training, as shown in Figure 5, we observe additional advantages. Even after accounting for the samples used during 412 pre-training, our method consistently outperforms the highly sample-efficient DreamerV3 baseline, 413 learning quickly and achieving high returns with minimal experience. This enhanced sample effi-414 ciency is observed across all tasks, highlighting the utility of the structured latent space to behavior 415 models. 416

417

395

396 397

398 399

Discovering Task-relevant Objects In Figure 6, we visualize an extract of the slot history via the 418 image reconstructions visualized in the first row. To visualize the attention pattern of the actor in 419 the current (rightmost) time-step we multiply the attention scores with the masks of the respective 420 objects and show them overlaid with the RGB reconstructions and as an individual colormap in the 421 second and third row respectively. This visualization shows the *Push-Specific* task and we find that 422 the model discovers task relevant objects automatically, ignoring information stored in slots that 423 represent distractor objects across all time-steps, while the robot and green cube receive the most 424 attention. While the recency bias induced by ALiBi is clearly visible, we find that the model learns to 425 pierce through it when necessary, attending to the target, which has been occluded for 15 time-steps 426 in the last time-step where it was visible. We see that the model is able to focus on task-relevant 427 information effectively, even when reasoning over long time sequences is required.

428

429 SAVi Finetuning One common limitation of prior works is that object-centric models are pre-430 trained on sequences with random behaviors but then frozen during training. This restricts their 431 applicability to tasks where the state distributions encountered by random and successful policies 432 are similar. With the *PickAndPlace* tasks, we explicitly violate this assumption since random be-



Figure 6: SOLD discovers objects relevant for task completion in an unsupervised manner over long horizons. We depict the normalized attention of the [out] token of the actor over the object tokens using Attention Rollout (Abnar & Zuidema, 2020). The full slot history is shown in Figure 15.

haviors are highly unlikely to pick and lift blocks off the table. Hence, SAVi has not seen blocks in the air, a configuration that will necessarily be reached by performant policies. Figure 7 illustrates the need to continually fine-tune the object-centric encoder-decoder model. While the fine-tuned SAVi model is able to reconstruct the lifted block accurately, the frozen variant fails to decode this configuration correctly. The target block effectively dissolves when lifted, hindering the discovery of behaviors associated with such states.

Generalization to Non-Object-Centric Environments While it is commonplace to evaluate 455 object-centric methods on environments and datasets that naturally lend themselves to such decom-456 positions, we aim to showcase the potential of our methods to generalize beyond this setting. To 457 this end, we tested our method on Meta-World to assess its performance on tasks involving objects 458 with complex shapes and colors, rather than simple uni-colored objects. Our method converges to 459 success rate of 100% on both the *Button-Press* and *Hammer* task. Further, we test on DM-Control 460 to evaluate generalization to problems that are vastly different and where reasoning over objects and interactions is less pronounced. On the Cartpole-Balance and Finger-Spin tasks, we reach returns of 461 497 and 645, respectively. We provide details about the environment decomposition and dynamics 462 prediction for all four tasks in Section F in the Appendix. 463

464 465

466

432

442 443

444

445

446 447

448

449

450

451

452

453 454

5 RELATED WORK

467 **Object-Centric Learning** In recent years, the field of unsupervised object-centric representation 468 learning from images and videos has gained significant attention (Yuan et al., 2023). Most existing 469 methods follow an encoder-decoder framework with a structured bottleneck composed of N latent 470 vectors called slots, where each of these slots binds to a different object in the input image. Slot-471 based methods have been widely applied for images (Burgess et al., 2019; Engelcke et al., 2020; Locatello et al., 2020; Singh et al., 2021; Engelcke et al., 2021; Singh et al., 2023; Biza et al., 2023) 472 and videos (Kipf et al., 2022; Singh et al., 2022; Elsayed et al., 2022; Bao et al., 2022; Zoran et al., 473 2021; Kabra et al., 2021; Creswell et al., 2021). However, despite their impressive performance 474 on synthetic datasets, they often fail to generalize to visually complex scenes. To overcome this 475 limitation, recent methods propose introducing some weak supervision (Elsayed et al., 2022; Bao 476 et al., 2023), levering large self-supervised pretrained encoders (Seitzer et al., 2023; Zadaianchuk 477 et al., 2024; Aydemir et al., 2023; Kakogeorgiou et al., 2024), or using diffusion models as slot 478 decoders (Jiang et al., 2023; Wu et al., 2023b; Singh et al., 2024).

479

Object-Centric Video Prediction Object-centric video prediction aims to understand the object dynamics in a video sequence with the goal of anticipating how these objects will move and interact with each other in future time steps. With this end, multiple methods propose to model and fore-cast the object dynamics using different architectures, including RNNs (Zoran et al., 2021; Nakano et al., 2023) transformers (Wu et al., 2023a; Villar-Corrales et al., 2023; Song et al., 2023; Gandhi et al., 2024; Daniel & Tamar, 2024; Nguyen et al., 2024; Petri et al., 2024) or state-space models (Jiang et al., 2024), achieving an impressive prediction performance on synthetic video datasets

491 492 493

494 495

496

497 498

499

500 501

502

503

504

505

506

507

508



Figure 7: Visualization of the full reconstruction and the slot that represents the target object for a frozen and finetuned SAVi model.

and learning representations that can help solve downstream tasks that require reasoning about objects properties and relationships (Wu et al., 2023a; Petri et al., 2024).

Model-based Reinforcement Learning Model-based reinforcement learning approaches aim to improve sample efficiency by learning environment dynamics. PlaNet (Hafner et al., 2019) introduced a latent dynamics model for efficient planning, while the Dreamer family (Hafner et al., 2020; 2021; 2023) incorporated this into an actor-critic framework. DreamerV2 and DreamerV3 introduced further improvements like categorical latent states and robustness techniques. DreamerV3 has shown superior performance in visual control tasks compared to model-free approaches, but uses holistic rather than object-centric state representations.

509 **Reinforcement Learning with Object-Centric Representations** Recent works have explored in-510 tegrating object-centric representations into RL frameworks. SMORL (Zadaianchuk et al., 2021) 511 and EIT (Haramati et al., 2024) combined object-centric representations with goal-conditioned 512 model-free RL for robotic manipulation. Yoon et al. (2023) investigated pre-training object-centric 513 representations for RL, showing benefits for relational reasoning tasks. The field of object-centric 514 model-based RL is still largely underexplored. One approach that can be categorized as such is FOCUS (Ferraro et al., 2023). However, unlike our method, FOCUS does not use the object-centric 515 states in forward prediction or action selection, but mainly for an exploration target. Further, FOCUS 516 requires supervision via ground-truth segmentation masks to learn the object-centric states. 517

518 519

520

6 CONCLUSION

521 We present SOLD, an object-centric model-based RL algorithm that learns directly from pixel inputs. By employing structured latent representations through slot-based dynamics models, our 522 method offers a compelling alternative to traditional, holistic approaches. While object-centric rep-523 resentations have been valued for their role in forward prediction (Villar-Corrales et al., 2023), we demonstrate their synergistic benefits in accelerating the learning of behavior models. SOLD 525 achieved strong performance across visual robotics environments, significantly outperforming the 526 state-of-the-art DreamerV3, particularly in tasks requiring relational reasoning. Additionally, the 527 learned behavior models exhibit interpretable attention patterns, explicitly focusing on task-relevant 528 parts of the visual scene.

529

530 **Limitations & Future Work** One limitation of our world model is that it generates predictions 531 in a deterministic manner. This can be a drawback in environments that are inherently stochastic 532 or highly unpredictable. We believe this is a key reason why SOLD outperforms DreamerV3 on 533 complex robotic manipulation tasks but struggles to match its performance on simpler tasks like 534 Cartpole-Balance, where minor variations in action sequences can lead to vastly different outcomes 535 over long horizons. Addressing this limitation by incorporating stochasticity into the prediction 536 model presents a promising direction for future work. A second limitation arises from the objectcentric encoder-decoder model we use. While SAVi performs well on the tasks we evaluated, scaling 537 it to complex real-world data remains a significant challenge. However, the core ideas of our method 538 are independent of the specific object-centric encoder-decoder model, and future work can easily 539 integrate more advanced models that address these scalability concerns.

540 REPRODUCIBILITY STATEMENT

With the goal of reproducible research, we conducted three runs per method and task, each with different random seeds to account for seed-dependent variability. Additionally, we thoroughly list in Appendix C the model architectures and hyper-parameters used in our experiments. Upon acceptance of the paper, we will open-source our code, environments and pretrained models.

547 548 REFERENCES

549

550

555

556

557

558

559

560 561

562

563 564

565

567

583

584

585

- Samira Abnar and Willem Zuidema. Quantifying attention flow in transformers. In Annual Meeting of the Association for Computational Linguistics (ACL), 2020.
- Eloi Alonso, Adam Jelley, 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.
 - Görkay Aydemir, Weidi Xie, and Fatma Guney. Self-supervised object-centric learning for videos. In Advances in Neural Information Processing Systems (NeurIPS), 2023.
 - Zhipeng Bao, Pavel Tokmakov, Allan Jabri, Yu-Xiong Wang, Adrien Gaidon, and Martial Hebert. Discovering objects that can move. In *Conference on Computer Vision and Pattern Recognition* (CVPR), 2022.
 - Zhipeng Bao, Pavel Tokmakov, Yu-Xiong Wang, Adrien Gaidon, and Martial Hebert. Object discovery from motion-guided tokens. In *Conference on Computer Vision and Pattern Recognition* (*CVPR*), 2023.
 - Peter W. Battaglia, Razvan Pascanu, Matthew Lai, Danilo Jimenez Rezende, and Koray Kavukcuoglu. Interaction networks for learning about objects, relations and physics. In Advances in Neural Information Processing Systems (NeurIPS), pp. 4502–4510, 2016.
- Erin L Beatty and Oshin Vartanian. The prospects of working memory training for improving
 deductive reasoning. *Frontiers in Human Neuroscience*, 9:56, 2015.
- Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement learning. In *International Conference on Machine Learning (ICML)*, 2017.
- Ondrej Biza, Sjoerd Van Steenkiste, Mehdi SM Sajjadi, Gamaleldin F Elsayed, Aravindh Mahen dran, and Thomas Kipf. Invariant slot attention: Object discovery with slot-centric reference
 frames. In *International Conference on Machine Learning (ICML)*, 2023.
- 576
 577 Christopher P Burgess, Loic Matthey, Nicholas Watters, Rishabh Kabra, Irina Higgins, Matt Botvinick, and Alexander Lerchner. MONet: Unsupervised scene decomposition and representation. *arXiv preprint arXiv:1901.11390*, 2019.
- Kyunghyun Cho. Learning phrase representations using RNN encoder-decoder for statistical
 machine translation. In *Conference on Empirical Methods in Natural Language Processing* (*EMNLP*), 2014.
 - Antonia Creswell, Rishabh Kabra, Chris Burgess, and Murray Shanahan. Unsupervised object-based transition models for 3D partially observable environments. In Advances in Neural Information Processing Systems (NeurIPS), 2021.
- Sebastian J Crutch, Sarah Connell, and Elizabeth K Warrington. The different representational frameworks underpinning abstract and concrete knowledge: Evidence from odd-one-out judgements. *Quarterly Journal of Experimental Psychology*, 62(7):1377–1390, 2009.
- Tal Daniel and Aviv Tamar. DDLP: Unsupervised Object-centric Video Prediction with Deep Dy namic Latent Particles. *Transactions on Machine Learning Research (TMLR)*, 2024.
- 593 Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers. In *International Conference on Learning Representations (ICLR)*, 2024.

594	Gamaleldin Elsayed, Aravindh Mahendran, Sjoerd Van Steenkiste, Klaus Greff, Michael C Mozer,	
595 596	and Thomas Kipf. SAVi++: Towards End-to-End Object-Centric Learning from Real-World	
597	videos. In Advances in Neural Information Processing Systems (NeurIPS), 2022.	
598	Martin Engelcke, Adam R Kosiorek, Oiwi Parker Jones, and Ingmar Posner. GENESIS: Gener-	
599	ative scene inference and sampling with object-centric latent representations. In International	
600	Conference on Learning Representations (ICLR), 2020.	
601	Martin Engelcke, Oiwi Parker Jones, and Ingmar Posner. GENESIS-V2: Inferring unordered ob-	
602	ject representations without iterative refinement. In Advances in Neural Information Processing	
603	Systems (NeurIPS), 2021.	
604	Stefano Ferraro Pietro Mazzaglia Tim Verbelen and Bart Dhoedt FOCUS: Object-centric world	
605	models for robotic manipulation. In Advances in Neural Information Processing Systems Work-	
606	shops (NeurIPSW), 2023.	
608	Sanket Gandhi, Samanyu Mahajan, Vishal Sharma, Rushil Gunta, Arnah Kumar Mondal, Parag	
609	Singla, et al. Learning disentangled representation in object-centric models for visual dynamics	
610	prediction via transformers. arXiv preprint arXiv:2407.03216, 2024.	
611		
612	David Ha and Jurgen Schmidhuber. Recurrent world models facilitate policy evolution. In Advances in Neural Information Processing Systems (NeurIPS) 2018	
613	in Neurai information 1 rocessing Systems (Neurin 5), 2018.	
614	Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James	
615	Davidson. Learning latent dynamics for planning from pixels. In International conference on	
616	Machine Learning (ICML), pp. 2555–2565. PMLR, 2019.	
617	Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learn-	
618	ing behaviors by latent imagination. In International Conference on Learning Representations	
619	(<i>ICLR</i>), 2020.	
621	Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering A	
622	discrete world models. In International Conference on Learning Representations (ICLR), 2021.	
623	Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains	
624	through world models. arXiv preprint arXiv:2301.04104, 2023.	
626	Nicklas Hansen, Xiaolong Wang, and Hao Su. Temporal difference learning for model predictive	
627	control. In International Conference on Machine Learning (ICML), 2022.	
628	Nicklas Hansen, Hao Su, and Xiaolong Wang. TD-MPC2: Scalable, robust world models for con-	
629	tinuous control. In International Conference on Learning Representations (ICLR), 2024.	
630	Dan Haramati, Tal Daniel, and Aviv Tamar. Entity-centric reinforcement learning for object manip-	
632	ulation from pixels. In International Conference on Learning Representations (ICLR), 2024.	
633	Jindong Jiang, Fei Deng, Gautam Singh, and Sungjin Ahn. Object-centric slot diffusion. In Ad-	
634	vances in Neural Information Processing Systems (NeurIPS), 2023.	
635	Jindong Jiang Fei Deng Gautam Singh Minseung Lee and Sungiin Ahn, SlotSSMs: Slot State	
636	Space Models. In Advances in Neural Information Processing Systems (NeurIPS), 2024.	
637		
630	Rishabh Kabra, Daniel Zoran, Goker Erdogan, Loic Matthey, Antonia Creswell, Matt Botvinick,	
640	iect representations via unsupervised video decomposition. In Advances in Neural Information	
641	Processing Systems (NeurIPS), 2021.	
642		
643	Ioannis Kakogeorgiou, Spyros Uidaris, Konstantinos Karantzalos, and Nikos Komodakis. SPOI: Self-training with patch-order permutation for object centric learning with autoregressive trans-	
644	formers. In Conference on Computer Vision and Pattern Recognition (CVPR), 2024	
645		
646	Elia Kaufmann, Leonard Bauersfeld, Antonio Loquercio, Matthias Müller, Vladlen Koltun, and	
647	Davide Scaramuzza. Champion-level drone racing using deep reinforcement learning. <i>Nature</i> , 620(7976):982–987, 2023.	

648 649	Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In <i>International Conference on Learning Representations (ICLR)</i> , 2015.
651 652 653	Thomas Kipf, Gamaleldin F Elsayed, Aravindh Mahendran, Austin Stone, Sara Sabour, Georg Heigold, Rico Jonschkowski, Alexey Dosovitskiy, and Klaus Greff. Conditional object-centric learning from video. In <i>International Conference on Learning Representations (ICLR)</i> , 2022.
654 655 656	Richard Li, Allan Jabri, Trevor Darrell, and Pulkit Agrawal. Towards practical multi-object ma- nipulation using relational reinforcement learning. In <i>International Conference on Robotics and</i> <i>Automation (ICRA)</i> , 2020.
657 658 659 660	Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 2020.
661 662	Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In Inter- national Conference on Learning Representations (ICLR), 2017.
663 664 665	Vincent Micheli, Eloi Alonso, and François Fleuret. Transformers are sample-efficient world models. In <i>International Conference on Learning Representations (ICLR)</i> , 2023.
666 667 668	Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Belle- mare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. <i>Nature</i> , 518(7540):529–533, 2015.
669 670 671	Akihiro Nakano, Masahiro Suzuki, and Yutaka Matsuo. Interaction-based disentanglement of en- tities for object-centric world models. In <i>International Conference on Learning Representations</i> (ICLR), 2023.
672 673 674 675	Trang Nguyen, Amin Mansouri, Kanika Madan, Khuong Duy Nguyen, Kartik Ahuja, Dianbo Liu, and Yoshua Bengio. Reusable slotwise mechanisms. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 2024.
676 677 678 679	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. PyTorch: An imperative style, high-performance deep learning library. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 2019.
680 681 682	Francesco Petri, Luigi Asprino, and Aldo Gangemi. Transformers and slot encoding for sample efficient physical world modelling. <i>arXiv preprint arXiv:2405.20180</i> , 2024.
683 684 685	Ofir Press, Noah Smith, and Mike Lewis. Train short, test long: Attention with linear biases enables input length extrapolation. In <i>International Conference on Learning Representations (ICLR)</i> , 2022.
686 687 688	Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering Atari, Go, chess and shogi by planning with a learned model. <i>Nature</i> , 588(7839):604–609, 2020.
690 691 692 693	Maximilian Seitzer, Max Horn, Andrii Zadaianchuk, Dominik Zietlow, Tianjun Xiao, Carl-Johann Simon-Gabriel, Tong He, Zheng Zhang, Bernhard Schölkopf, Thomas Brox, et al. Bridging the gap to real-world object-centric learning. In <i>International Conference on Learning Representations (ICLR)</i> , 2023.
694 695 696	Gaurav Sharma, Wencheng Wu, and Edul N Dalal. The CIEDE2000 color-difference formula: Implementation notes, supplementary test data, and mathematical observations. <i>Color Research & Application</i> , 30(1):21–30, 2005.
697 698 699	David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of Go with deep neural networks and tree search. <i>Nature</i> , 529(7587):484–489, 2016.
700	Gautam Singh, Fei Deng, and Sungjin Ahn. Illiterate dall-e learns to compose. In International Conference on Learning Representations (ICLR), 2021.

702 703 704	Gautam Singh, Yi-Fu Wu, and Sungjin Ahn. Simple unsupervised object-centric learning for com- plex and naturalistic videos. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 2022.
705 706 707	Gautam Singh, Yeongbin Kim, and Sungjin Ahn. Neural systematic binder. In International Con- ference on Learning Representations (ICLR), 2023.
708 709	Krishnakant Singh, Simone Schaub-Meyer, and Stefan Roth. Guided latent slot diffusion for object- centric learning. <i>arXiv preprint arXiv:2407.17929</i> , 2024.
710 711 712 713	Yeon-Ji Song, Hyunseo Kim, Suhyung Choi, Jin-Hwa Kim, and Byoung-Tak Zhang. Learning object motion and appearance dynamics with object-centric representations. In Advances in Neural Information Processing Systems Workshops (NeurIPSW), 2023.
714	Elizabeth S Spelke. Principles of object perception. Cognitive Science, 14(1):29-56, 1990.
715 716 717	Richard S Sutton and Andrew G. Barto. <i>Reinforcement learning: An introduction</i> . A Bradford Book, 2018.
718 719 720 721	Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Bud- den, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. DeepMind control suite. arXiv preprint arXiv:1801.00690, 2018.
722 723	Emanuel Todorov, Tom Erez, and Yuval Tassa. MuJoCo: A physics engine for model-based control. In International Conference on Intelligent Robots and Systems (IROS), 2012.
724 725 726	Dani Valevski, Yaniv Leviathan, Moab Arar, and Shlomi Fruchter. Diffusion models are real-time game engines. <i>arXiv preprint arXiv:2408.14837</i> , 2024.
727 728 729	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in Neural Information Processing Systems (NeurIPS), 2017.
730 731 732	Angel Villar-Corrales, Ismail Wahdan, and Sven Behnke. Object-centric video prediction via de- coupling of object dynamics and interactions. In <i>International Conference on Image Processing</i> (<i>ICIP</i>), 2023.
733 734 735	J Beau W Webber. A bi-symmetric log transformation for wide-range data. <i>Measurement Science and Technology</i> , 24(2):027001, 2012.
736 737 738 720	Ziyi Wu, Nikita Dvornik, Klaus Greff, Thomas Kipf, and Animesh Garg. Slotformer: Unsupervised visual dynamics simulation with object-centric models. In <i>International Conference on Learning Representations (ICLR)</i> , 2023a.
740 741 742	Ziyi Wu, Jingyu Hu, Wuyue Lu, Igor Gilitschenski, and Animesh Garg. Slotdiffusion: Object- centric generative modeling with diffusion models. In <i>Advances in Neural Information Processing</i> <i>Systems (NeurIPS)</i> , 2023b.
743 744 745	Jaesik Yoon, Yi-Fu Wu, Heechul Bae, and Sungjin Ahn. An investigation into pre-training object- centric representations for reinforcement learning. In <i>International Conference on Machine</i> <i>Learning (ICML)</i> , 2023.
747 748 749	Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-World: A benchmark and evaluation for multi-task and meta reinforcement learning. In <i>Conference on Robot Learning (CoRL)</i> , 2019.
750 751 752	Jinyang Yuan, Tonglin Chen, Bin Li, and Xiangyang Xue. Compositional scene representation learning via reconstruction: A survey. <i>Transactions on Pattern Analysis and Machine Intelligence</i> (<i>TPAMI</i>), 45(10):11540–11560, 2023.
753 754 755	Andrii Zadaianchuk, Maximilian Seitzer, and Georg Martius. Self-supervised visual reinforcement learning with object-centric representations. In <i>International Conference on Learning Representations (ICLR)</i> , 2021.

756 757 758 750	Andrii Zadaianchuk, Maximilian Seitzer, and Georg Martius. Object-centric learning for real-world videos by predicting temporal feature similarities. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> , 2024.
759 760 761	Biao Zhang and Rico Sennrich. Root mean square layer normalization. In Advances in Neural Information Processing Systems (NeurIPS), 2019.
762	Daniel Zoran, Rishabh Kabra, Alexander Lerchner, and Danilo J Rezende. PARTS: Unsupervised
763	segmentation with slots, attention and independence maximization. In International Conference
764	on Computer Vision (ICCV), 2021.
765	
766	
767	
768	
769	
770	
771	
772	
773	
774	
776	
777	
778	
779	
780	
781	
782	
783	
784	
785	
786	
787	
788	
789	
790	
791	
792	
793	
794	
795	
796	
797	
790	
800	
801	
802	
803	
804	
805	
806	
807	
808	
809	

810	A NOTATION	
811		
812		Slot Attention for Video (SAVI)
813	D_Z	The slot dimension
815	N	The number of slots
816	$oldsymbol{Z}_t$	A set of slots $\boldsymbol{Z}_t \in \mathbb{R}^{N imes D_Z}$ at time-step t
81 <i>1</i> 818	$oldsymbol{Z}_{0:t}$	A history of slot-sets up to time-step t
819	e_ψ	A SAVi encoder that maps \boldsymbol{o}_t to \boldsymbol{Z}_t
820	d_ψ	A SAVi decoder that reconstructs \boldsymbol{o}_t from \boldsymbol{Z}_t
822	$oldsymbol{F}_t$	Features obtained by encoding images
823	L	Number of spatial locations in F
824		
		D!
825		Reinforcement Learning
825 826 827	S	Reinforcement Learning The number of seed frames
825 826 827 828	$S \ T$	Reinforcement Learning The number of seed frames The imagination horizon
825 826 827 828 829	$S \\ T \\ o_t$	Reinforcement Learning The number of seed frames The imagination horizon An image observation
825 826 827 828 829 830 831	$S \\ T \\ o_t \\ a_t$	Reinforcement Learning The number of seed frames The imagination horizon An image observation An action command
825 826 827 828 829 830 831 832	$S \\ T \\ o_t \\ a_t \\ r_t$	Reinforcement Learning The number of seed frames The imagination horizon An image observation An action command A reward
825 826 827 828 829 830 831 832 833	$egin{array}{ccc} S & & \ T & & \ o_t & & \ a_t & & \ r_t & & \ \gamma & & \end{array}$	Reinforcement Learning The number of seed frames The imagination horizon An image observation An action command A reward A scalar discount factor
825 826 827 828 829 830 831 832 833 833	$egin{array}{cccc} S & T & & \ T & oldsymbol{o}_t & & \ oldsymbol{a}_t & & \ oldsymbol{r}_t & & \ oldsym$	Reinforcement Learning The number of seed frames The imagination horizon An image observation An action command A reward A scalar discount factor
825 826 827 828 830 831 832 833 834 835	$egin{array}{ccc} S & T & & \ T & oldsymbol{o}_t & & \ oldsymbol{a}_t & & \ r_t & & \ \gamma & & \ H & & \ \end{array}$	Reinforcement Learning The number of seed frames The imagination horizon An image observation An action command A reward A scalar discount factor The entropy of a probability distribution
825 826 827 828 830 831 832 833 834 835 836	$egin{array}{c} S \ T \ m{o}_t \ m{a}_t \ r_t \ \gamma \ m{H} \ f^{ ext{MLP}}_lpha \end{array}$	Reinforcement LearningThe number of seed framesThe imagination horizonAn image observationAn action commandA rewardA scalar discount factorThe entropy of a probability distributionAn MLP head that belongs to parameter group α
825 826 827 828 830 831 832 833 834 835 836 837	$egin{array}{ccc} S & T & & \ m{v}_t & & \ m{a}_t & & \ m{r}_t & & \ m{\gamma} & & \ m{H} & & \ f_{lpha}^{ ext{MLP}} & & \ m{h}_t & & \ m{h}_t \end{array}$	Reinforcement Learning The number of seed frames The imagination horizon An image observation An action command A reward A scalar discount factor The entropy of a probability distribution An MLP head that belongs to parameter group α A processed output token

910

840 841 842

EXTENDED RELATED WORK Β

843 **Model-based Reinforcement Learning** Model-based methods hold the potential to significantly improve the sample efficiency of RL algorithms, and recent years have seen several key contributions 844 advancing this area. 845

846 Pioneering work by Ha & Schmidhuber (2018) introduced the concept of a recurrent generative 847 model, termed a world model, which captures the dynamics of visual RL environments. By encoding 848 high-dimensional observations into a compact latent representation, this model enables RL agents 849 to train policies entirely within imagined rollouts.

850 The Planning Network (PlaNet) (Hafner et al., 2019) introduced a recurrent state-space model 851 (RSSM) that predicts future states directly in a compact latent space, avoiding the costly step of 852 decoding full observations. This architecture enables efficient planning of action sequences but 853 is limited by short planning horizons. Building on this, Dreamer (Hafner et al., 2020) integrates 854 planning and learning by training agents within a learned world model, overcoming PlaNet's shortsighted horizon. Subsequent versions, DreamerV2 (Hafner et al., 2021) and DreamerV3 (Hafner 855 et al., 2023), improved robustness and generalization through enhanced representation learning and 856 optimization techniques, achieving state-of-the-art results across diverse RL benchmarks. 857

858 Temporal Difference Learning for Model Predictive Control (TD-MPC) (Hansen et al., 2022) intro-859 duced a task-oriented latent dynamics model to optimize trajectories directly within the latent space 860 of a world model. Unlike earlier approaches, TD-MPC avoids reconstructing full observations, 861 instead focusing the world model on reward-predictive elements through a loss applied to reward and value predictions. TD-MPC2 (Hansen et al., 2024) builds on this by introducing scalability 862 improvements, enabling superior performance with larger model sizes and demonstrating a single 863 agent's ability to generalize across multiple tasks and action spaces.

865	(a) SAVi		(b) Object-centric dynamics		(c) Slot Aggregation Transformer	
867	Hyper-Param.	Value	Hyper-Param.	Value	Hyper-Param.	Value
868	Slot Dim. D_Z	128	# Layers	4	# Layers	4
869	# Slots N	2-10	# Heads	8	# Heads	8
870	Slot Init.	Learned	Token Dim.	256	Token Dim.	256
871	# Iters.	3/1	MLP Dim.	512	MLP Dim.	512

Table 1: Implementation details for each of SOLD modules.

Inspired by the success of Transformers on sequence modeling tasks, Micheli et al. (2023) proposed 874 IRIS, a method combining a discrete autoencoder with an autoregressive Transformer to model 875 environment dynamics. The autoencoder tokenizes images into a discrete set of representations, 876 while the Transformer learns temporal dynamics across these tokens. IRIS demonstrated visually 877 and temporally accurate predictions of game dynamics in Atari environments. While IRIS shares 878 similarities with our approach - encoding an image into a set-based representation and predicting it 879 forward using a Transformer – it lacks the object-centric interpretability afforded by our model. 880

Recently, Alonso et al. (2024) proposed diffusion as a promising alternative to discretization of the latent space of the world model. DIAMOND uses a diffusion-based conditional generative model, $p(x_{t+1} \mid x_{\leq t}, a_{\leq t})$, to produce visually precise next-frame predictions. The authors demonstrate the potential of their world model to simulate complex 3D environments by learning a realistic game-engine from static Counter-Strike: Global Offensive gameplay. Valevski et al. (2024) also propose to use diffusion to create high-quality visual predictions. Specifically, they demonstrate the potential of diffusion models to serve as real-time game engines.

С IMPLEMENTATION DETAILS 889

891 In this section, we describe the network architecture and training details for each of the SOLD com-892 ponents. Our models are implemented in PyTorch (Paszke et al., 2019), have 12 million learnable parameters, and are trained on a single NVIDIA A-100 GPU with 40GB of VRAM. A summary of 893 the model implementation details is listed in Table 1. 894

C.1 SLOT ATTENTION FOR VIDEO 896

We closely follow Kipf et al. (2022) for the implementation of the Slot Attention for Video 898 (SAVi) decomposition model, including their proposed CNN-based encoder e_{ub} and decoder d_{ub} . 899 the transformer-based predictor and the Slot Attention corrector. We employ between 2 and 10 900 (depending on the environment) 128-dimensional object slots, whose initialization is learned via 901 backpropagation. We empirically verified that learning the initial slots performs more stable than 902 the usual random initialization. Furthermore, we use three Slot Attention iterations for the first video 903 frame in order to obtain a good initial decomposition, and a single iteration for subsequent frames, 904 which is enough to update the slot state given the observed features. 905

906 C.2 **OBJECT-CENTRIC DYNAMICS MODEL** 907

908 Our object-centric dynamics model is based on the OCVP-Seq (Villar-Corrales et al., 2023) archi-909 tecture, which is a transformer encoder employing sequential and relational attention mechanisms 910 in order to decouple the processing of temporal dynamics and interactions, and has been shown to achieve interpretable and temporally consistent predictions. We use 4 transformer layers employing 911 256-dimensional tokens, 8 attention heads, and using a hidden dimension of 512 in the feed-forward 912 layers. 913

914

916

864

881

882

883

884

885

886

887 888

890

895

897

915 SLOT AGGREGATION TRANSFORMER C.3

The Slot Aggregation Transformer (SAT) forms the architectural backbone for the reward, value and 917 action models. This module aggregates information from object slots across multiple time steps to



Figure 8: The *Slot Aggregation Transformer* applies causal masking and ensures that output and
register token do not attend to themselves on other time-steps. The recency bias induced by ALiBi is
visualized through the color gradient in the attention mask, with lighter shades of blue corresponding
to a higher negative bias on the attention scores.

produce output tokens that are subsequently fed to MLP heads in order to predict rewards, values, or actions. An overview of our SAT module is depicted in Figure 8.

SAT is a causal transformer encoder module that receives as input a history of object slots, as well as a learnable output token [out] for each time step, which is responsible for producing the final output for the corresponding time step. Additionally, we append to the SAT inputs a number of register tokens [reg] per time-step, which have been shown to aid with processing in attention-based models by offloading intermediate computations from the output tokens and helping the module focus on relevant slots (Darcet et al., 2024).

To encode the positional information into SAT, we employ *Attention with Linear Biases* (Press et al., 2022) (ALiBi), which introduces linear biases directly into the attention scores, effectively encoding token recency. This approach helps the model deal with sequences of varying length, as well as generalize to longer sequences than those seen during training, thus outperforming absolute positional encodings.

For our experiments, SAT is implemented with 4 transformer encoder layers with causal selfattention, RMS-Normalization layers (Zhang & Sennrich, 2019), 8 attention heads, a token dimension of 256, and a hidden dimensionality in the feed-forward layers of 512. We set the number of learnable register token per time step to 4. Furthermore, we enforce in our causal attention masks that tokens belonging to time step t cannot directly interact with previous output and register tokens.

C.4 TRAINING DETAILS

SAVi Pretraining SAVi is pretrained for object-centric decomposition on approximately one million frames for 400,000 gradient steps. We use the Adam optimizer (Kingma & Ba, 2015), a batch size of 64 and a base learning rate of 10^{-4} , which is first linearly warmed-up during the first 2,500 training steps, followed by cosine annealing (Loshchilov & Hutter, 2017) for the remaining of the training procedure. We perform gradient clipping with a maximum norm of 0.05. 972 SOLD Training SOLD is trained using the Adam optimizer (Kingma & Ba, 2015) and different 973 learning rates for each component: 10^{-4} for the dynamics and rewards models, and $3 \cdot 10^{-5}$ for 974 training the action and value models, as well as for fine-tuning the SAVi encoder. To stabilize 975 training, we perform gradient clipping with maximum norm of 0.05 for the SAVi model, 3.0 for 976 the transition model, and 10.0 for the reward, value, and action models. For all components, we also use learning rate warmup for the first 2,500 gradient steps. Additionally, we implement the 977 exponential moving average (EMA) for the target value network with a decay rate of 0.98. We use 978 an imagination horizon of 15 steps for behavior learning, and the λ -parameter is set to 0.95. 979

980 981

982

985

D BASELINES

In our experiments we compare our approach with two different baseline models, namely the SoTA
 model-based RL baseline *DreamerV3* and a *Non-Object-Centric* variant of our proposed model:

986 **DreamerV3** DreamerV3 (Hafner et al., 2023) is a SoTA model-based reinforcement learning algorithm that learns behaviors from visual inputs without requiring task-specific inductive biases or 987 extensive environment interaction. It builds a world model that predicts future states and rewards, 988 which is then used to simulate potential outcomes and guide the agent's decision. DreamerV3 989 leverages latent dynamics and a compact, holistic representation of the environment for an efficient 990 exploration, while showing desirable properties such as sample efficiency, scalability, and gener-991 alization across a wide range of complex tasks. We select the 12 million parameter variant so as 992 to match the parameter count of our proposed model. For further details, we refer to Hafner et al. 993 (2023). 994

995 **Non-Object-Centric Baseline** This baseline model follows the same general framework as our 996 proposed model, but replaces the object-centric SAVi encoder and decoder with a simple convo-997 lutional auto-encoder while keeping the remaining modules unchanged; thus allowing us to ablate 998 the effect of object-centric representations for model-based reinforcement learning. The CNN auto-999 encoder used in this baseline consists of an encoder comprised of four strided convolutional layers 1000 with 64, 128, 256, and 512 channels respectively, each followed by batch normalization and a ReLU. The output of the final convolutional layer is flattened and fed through a linear layer to produce a 1001 512-dimensional latent vector. The decoder mirrors the encoder structure, reconstructing the obser-1002 vations from the latent representation through the use of four transposed convolutional layers. To 1003 compensate for the lack of multiple latent vectors and to ensure a fair comparison, we increase the 1004 capacity of this baseline model by scaling the actor, critic, and dynamics models. Specifically, we 1005 increase the token dimension from 256 to 512, as well as the MLP hidden dimension from 512 to 1006 1024. The total parameter count for this baseline is approximately 60 million, thus being five times 1007 larger than our proposed method and the DreamerV3 baseline.

1008 1009 1010

E ENVIRONMENTS

In this section, we provide further details about our proposed suite of environments, which includes eight object-centric robotic control tasks designed to test relational reasoning and manipulation capabilities. The environments, which are inspired by Li et al. (2020) and are simulated using Mu-JoCo (Todorov et al., 2012), follow the same basic structure, consisting of a robot arm mounted on a base, positioned near a table where the manipulation tasks take place.

In all environments, the robot is controlled by a 4-dimensional action vector $a = \begin{bmatrix} a_x, a_y, a_z, a_{grip} \end{bmatrix} \in [-1, 1]^4$, where the first three components represent the desired movement direction of the end-effector, whereas the fourth component controls the gripper. On the *Reach* and *Push* tasks, commands to the gripper are ignored, with the gripper fixed in a closed configuration, as gripping is not required to solve these tasks.

- 1022 For all tasks, we define the following constants:
- 1023 1024

- $t_1 = 20$ and $t_2 = 10$: Temperature parameters that determine the steepness of the reward function.
 - $d_m = 0.05$: Distance threshold (in meters) for considering a task successful.



Figure 9: Object-centric SAVi decomposition of a video frame. We show the masked RGB image and the segmentation mask corresponding to each object slot. The masked RGB images are combined to reconstruct the observed frame.

Reach In *Reach* tasks, the agent must identify a spherical target among several distractors and mode the end-effector to its location. The reward is calculated as:

$$r = \exp(-t_1 \cdot ||\boldsymbol{p}_e - \boldsymbol{p}_t||_2), \tag{12}$$

where p_e is the position of the end-effector and p_t is the target position. The success at the end of an episode is defined as:

$$\operatorname{success} = \begin{cases} 1 & \text{if } ||\boldsymbol{p}_e - \boldsymbol{p}_t||_2 < d_m \\ 0 & \text{otherwise} \end{cases}$$
(13)

Push & PickAndPlace Both *Push* and *PickAndPlace* correspond to reasoning and manipulation tasks where the agent must identify a single block among several distractors and move it to a target location. In *Push* tasks, the agent can slide the block to the target position on the table without using its gripper; whereas in *PickAndPlace* the target location can be above the table, thus needing to grasp the block in order to lift it to the target position. In both task variants the reward is calculated as:

$$r = 0.9 \cdot \exp(-t_1 \cdot ||\boldsymbol{p}_c - \boldsymbol{p}_t||_2) + 0.1 \cdot \exp(-t_2 \cdot ||\boldsymbol{p}_e - \boldsymbol{p}_c||_2), \tag{14}$$

where p_e is the position of the end-effector, p_t is the target position, and p_c is the block position. The success at the end of an episode is defined as:

$$\operatorname{success} = \begin{cases} 1 & \text{if } ||\boldsymbol{p}_c - \boldsymbol{p}_t||_2 < d_m \\ 0 & \text{otherwise} \end{cases}$$
(15)

1060 1061

1063

1065

1054 1055

1056

1057 1058 1059

1035

1036

1037 1038 1039

1040

1041 1042 1043

1044

1062 F ADDITIONAL RESULTS

1064 F.1 OBJECT-CENTRIC DECOMPOSITION

Figure 9 depicts the object-centric decomposition of a video frame obtained by SAVi. SAVi parses the input frame into per-object RGB reconstructions and alpha masks, which can be combined via a weighted sum in order to accurately reconstruct the observed video frame. Notably, SAVi assigns an object slot to the scene background, five slots to different blocks, one slot to the red target, and one slot to the robot arm. The sharp object masks demonstrate that SAVi isolates object-specific information in each slot, which is beneficial for downstream applications such as behavior learning, allowing the agent to reason about object properties and their relationships while abstracting taskirrelevant details.

1073

1074 F.2 OPEN-LOOP PREDICTION 1075

We visualize action-conditional open-loop predictions on the *Push Specific* (Figure 10), *Button-Press* (Figure 11), *Hammer* (Figure 12), *Cartpole-Balance* (Figure 13), and *Finger-Spin* (Figure 14)
environments. More precisely, we depict the ground truth sequence, the predicted video frames, the predicted instance segmentation of the scene, obtained by assigning a distinct color to each object mask, and the object reconstruction for each slot.



Figure 10: Open-loop prediction on the *Push-Specific* task. We visualize the ground truth, predicted
frames, segmentation obtained by assigning different colors to each object mask, and per-object
reconstructions. In this sequence, SOLD assigns one slot to the background, one slot for the robot,
one slot for the target, and four different slots for blocks, while one slot remains empty.

In all examples, our model parses the scene into sharp and accurate object representations, and
models the action-conditional object dynamics and interactions in order to accurately predict the
future video frames while preserving object-centric representations. We emphasize SOLD's ability
to capture complex physical interactions, such as pushing a block to a target location (Fig. 10),
pressing a button (Fig. 11) or hitting a nail with a hammer (Fig. 12). Furthermore, we showcase that
SOLD can generalize to diverse non-object-centric environments (Fig. 13 and Fig. 14).



Figure 11: Open-loop prediction and decomposition results on the *Button-Press* task. We visualize the ground truth and predicted video frames, instance segmentation obtained by assigning a different color to each object mask, and per-object reconstructions. SOLD assigns a slot for the scene background, two slots for different robot parts, and a slot for the button-box.



Figure 12: Open-loop prediction and decomposition results on the *Hammer* task. We visualize the ground truth and predicted video frames, instance segmentation obtained by assigning a different color to each object mask, and per-object reconstructions. SOLD assigns a slot for the scene background, three slots for different robot parts, a slot for the hammer, and a slot for the nail-box.

1265

1266

1267



Figure 13: Open-loop prediction and decomposition results on the *Cartpole-Balance* task. We visualize the ground truth and predicted video frames, instance segmentation obtained by assigning a different color to each object mask, and per-object reconstructions. SOLD assigns a slot for the scene background and a slot for the cart-pole. Notably, the slot represents the object along with its shadow.



Figure 14: Open-loop prediction and decomposition results on the *Finger-Spin* task. We visualize the ground truth and predicted video frames, instance segmentation obtained by assigning a different color to each object mask, and per-object reconstructions. SOLD assigns a slot for the scene background, a slot for the finger, and a slot for the spinning target. Notably, the slots represent the objects along with their corresponding shadows.

1	296
1	297
1	298
1	299
1	300
1	301
1	302
1	303
1	304
1	305
1	306
1	307
1	308
1	309
1	310
1	311
1	312
1	313
1	314
1	315
1	316
1	317
1	318
1	319
1	320
1	321
1	322
1	323
1	324
1	325
1	326
1	327
1	328



Figure 15: Original, uncut slot-history used in Figure 6. The full rollout highlights both the recency bias introduced by ALiBi and the model's ability to overcome this bias when task-relevant information must be retrieved from slots far in the past.



Figure 16: Visualizing the actor's attention over the slot-history on the *PickAndPlace-Specific* task reveals that the robot's gripper and target block in the current time-step receive the most attention, highlighting the model's ability to focus on task-relevant components. Further, the red target sphere is mostly occluded in the current time-step making it difficult to reconstruct its position accurately. However, the model learns to integrate information from previous steps, where it was still clearly visible.



Figure 17: Training curves for the non-object-centric environments studied in this work.

1352 F.3 QUANTITATIVE EVALUATION

1350 1351

1360

1361 1362

1363

1364

1365

1378

We present the training performance of SOLD on non-object-centric environments for a single seed
in Figure 17. Although SOLD is not expected to outperform baseline methods in environments that
do not rely on relational reasoning or significantly benefit from object-centric decomposition, our
results highlight the potential of applying object-centric representations to various visual domains.
This is underscored by SOLD's ability to generalize to the complex hammer and button press tasks,
converging to a success rate of 100% on both problems.

F.4 GENERALIZATION TO DIFFERENT COLORS

To investigate the generalization beyond the training distribution, we evaluate the performance of all methods on a different color set for all tasks. The results are shown in Figure 18. When comparing to the original evaluation in Figure 4 we observe a very small drop in performance for both SOLD and DreamerV3, indicating robustness to changes in the object color.



1379 F.5 SAVI FINE-TUNING 1380

Figure 19 illustrates the importance of fine-tuning the object-centric encoder-decoder model with
 another example. Without fine-tuning, the blue color, which appears similarly on both colored blocks
 and the robot arm, leads to an even more drastic degradation of the reconstructions, where the robot



Figure 19: Comparison of fine-tuned and frozen SAVi models on PickAndPlace-Distinct. We visualize the full reconstruction and the slot that reconstructs the cube that is being lifted for both models. When the blue block is lifted off the table, the frozen model merges it with blue elements from the robot arm, deteriorating the prediction and hallucinating the arm going between the gripper fingers. The fine-tuned model, on the other hand, is able to reconstruct the sequence accurately.

itself is no longer accurately captured. In contrast, the fine-tuned model is able to reconstruct the sequence accurately.