

Representing Positional Information in Generative World Models for Object Manipulation

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** The ability to predict outcomes of interactions between embodied
2 agents and objects is paramount in the robotic setting. While model-based control
3 methods have started to be employed for tackling manipulation tasks, they have
4 faced challenges in accurately manipulating objects. As we analyze the causes
5 of this limitation, we identify the cause of underperformance in the way current
6 world models represent crucial positional information, especially about the target’s
7 goal specification for object positioning tasks. We propose two solutions for
8 generative world models: position-conditioned (PCP) and latent-conditioned (LCP)
9 policy learning. In particular, LCP employs object-centric latent representations
10 that explicitly capture object positional information for goal specification. This
11 naturally leads to the emergence of multimodal capabilities.

12 1 Introduction

13 Among RL algorithms, model-based approaches aim to provide greater data efficiency compared
14 to their model-free counterparts [1, 2]. With the advent of world models (WM) [3], model-
15 based agents have demonstrated impressive performance across various domains [4–7], including
16 real-world robotic applications [8, 9].

17 When considering robotic object manipulation
18 tasks, it seems natural to consider an object-
19 centric approach to world modeling. Object-
20 centric world models, like FOCUS [10] learn
21 a distinct dynamical latent representation per
22 **object**. This contrasts with the popular Dreamer
23 method [6], where a single **flat** representation,
24 referring to the whole scene is extracted.

25 Model-based generative agents, like Dreamer
26 and FOCUS, learn a latent model of the envi-
27 ronment dynamics by reconstructing the agent’s
28 observations and use it to generate latent se-
29 quences for learning a behavior policy in imagi-
30 nation [4, 11, 6]. However, these kinds of agents
31 have shown consistent issues in succeeding in
32 object manipulation tasks, both from proprio-
33 ceptive/vector inputs [12] and from images [13].

34 After analyzing the causes of failure of generative agents, we propose two solutions to improve
35 performance:

- 36 • a simpler solution, where the target is expressed as a vector of spatial coordinates, that
37 presents no major changes to the model architecture and minimal changes to policy learning;

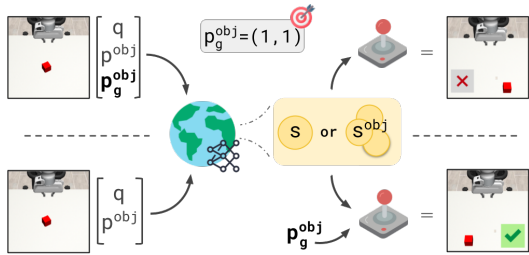


Figure 1: The world model compresses input observations into a single or per object latent state representation. The compressed representation serves as input to the policy for action selection. **(top)** Goal information is provided through the input state vector. **(bottom):** Both single and object-centric representations can be paired to a **target-conditioned policy**.

- a tailored solution employing an object-centric approach that integrates positional information about the objects into the latent space of the world model. This approach enables the possibility to specify goals through multimodal targets, e.g. vector inputs or visual goals.

2 Analysis of the Current Limitations

To provide insights into the limitations of current world model-based agents in object-positioning tasks, we consider the performance of Dreamer and FOCUS on a pose-reaching and an object-positioning task. For pose-reaching, we opted for the Reacher environment from the DMC suite [14]. In this task, we consider the end-effector of the manipulator as the entity to be positioned at the target location. For the more complex object positioning task, we opted for a cube-manipulation task from Robosuite [15]. The given cube has to be placed at the specified target location to succeed in the task.

In both environments, the target position is uniformly sampled within the workspace at every new episode. We test the environments in two different scenarios: first, with a virtual visual target that is rendered in the environment, and second, without a visual target, where the target location is provided only as a vector in the agent’s inputs. Training details are provided in appendix G. Based on Fig. 2, we highlight the significant gap in performance between the tasks with the virtual visual targets rendered in the environment and the tasks using only spatial coordinates as a target. The agents struggle to solve the tasks without a visual target. It can also be noticed a negative correlation between the agents’ ability to reconstruct positional information and the performance on the task.

There is a significant difference in the relative significance of the target information compared to the entire observation, in terms of their dimensionality. The information pertaining to a positional target comprises a maximum of three values (i.e., the xyz coordinates of the target). Conversely, when considering a visual cue, there are three values (i.e., RGB values) for each pixel that represents the target cue. Consequently, the relative significance of the target information is, at least, greater in the case of a large visual target, i.e. larger than a single pixel. This difference in the dimensionality affects the computation of the loss, and thus the weight of each component in the decoder’s loss. For the entity, the agents have access to this information in the visual observation. Indeed, it’s not surprising that both agents reconstruct the entity position accurately. To confirm our hypothesis that the improved predictions are due to the greater significance of the visual targets in the overall loss, we provide additional experiments in appendix C.

Discussion. A concurrent work [16] conducted an extensive study between the interplay of the reward and the observation loss in a world model. Our analysis provides an additional insight, as we identify within the observation loss, an unbalance between the different decoded components. In this work, rather than focussing on how to balance the losses (Appendix D), we consider different approaches to alleviate this issue. The central idea is to find alternative ways to provide positional information about the target directly to the reward computation and policy learning modules, rather than relying on the reconstruction of the targets obtained by the model.

3 Conditioned Policy

Position Conditioned Policy (PCP). The first declination of our proposed solutions is the conditioning of the policy directly on the positional coordinates of the desired target. By default, the world

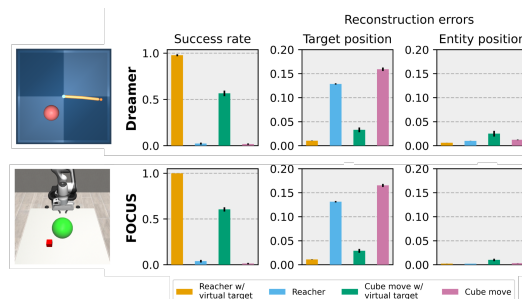


Figure 2: **(left)**: examples of virtual targets. **(top-right)**: Dreamer’s success rate and reconstruction performance over target and entity position (end-effector position for reacher and cube position for the cube environment). **(bottom-right)**: Equivalent for the FOCUS object-centric model.

86 model encodes the target’s positional information in the latent states, which are then fed to the policy
 87 for behavior learning. Instead, as shown in the bottom of Fig. 1, we propose to concatenate the object
 88 positional coordinates p_g^{obj} to the latent states s_t as an input to the policy network. We refer to this
 89 strategy as *Position-Conditioned Policy (PCP)*: $\pi_{PCP}(a_t|s_t, p_g^{obj})$

90 When employing PCP, the policy has direct access to the target’s positional information p_g^{obj} . This
 91 can also be leveraged for reward computation. Rather than learning a reward head, we can use the
 92 world model’s decoder to predict the object’s position at time t, obtaining \hat{p}_t^{obj} . Then, the reward
 93 r_{PCP} can be estimated as the distance between the target given to the policy and the reconstructed
 94 position of the entity of interest: $r_{PCP} = dist(\hat{p}_t^{obj} - p_g^{obj})$

95 **Latent Conditioned Policy (LCP).** Condition-
 96 ing the policy on explicit features has its limita-
 97 tions, particularly when extending features
 98 beyond positional ones, or when working with
 99 different goal specifications, e.g. visual ones.
 100 Therefore, expressing features implicitly could
 101 represent a more robust approach. To address
 102 this, we propose a latent conditioned method
 103 for behavior learning. This approach is anal-
 104 ogous to the one adopted in LEXA [17] for
 105 goal-conditioned behavior learning. However,
 106 we tailor our strategy for object manipulation
 107 by designing an object-centric approach. We
 108 refer to our novel implementation as *Latent-*
 109 *Conditioned Policy (LCP)*.

110 In LEXA, policy conditioning occurs on the
 111 entire (flat) latent state, using either cosine or
 112 temporal distance methods. However, in manipulation tasks involving small objects, the cosine
 113 approach is inadequate because it prioritizes matching the robot’s position over visually smaller
 114 aspects of the scene, such as an object’s position, rather than on bigger visual components of the
 115 scene, e.g. the robot pose. The temporal approach was introduced to mitigate this issue. However,
 116 this approach generally requires a larger amount of data to converge, as the training signal is less
 117 informative, being based only on the temporal distance from the goal [17]. We argue that object-
 118 centric latent representations offer greater flexibility to condition the policy, thanks to the disentangled
 119 latent information. With LCP, we can condition the policy solely on the object’s latent states, enabling
 120 fine-grained target conditioning focused exclusively on the entity of interest.

121 **Latent Positional Encoder.** To obtain object latent features for a given target position, we introduce
 122 the Latent Positional Encoder model, as shown in Fig. 3. This model enables inferring an object’s
 123 latent state directly from the object’s positional information, namely $p(\hat{s}_t^{obj}|p_t^{obj})$.

124 During training, the latent positional encoder is trained to minimize the negative cosine distance
 125 between the predicted and the reference object latent state: $\mathcal{L}_{pos} = -\frac{\hat{s}_t^{obj} \cdot s_t^{obj}}{\|\hat{s}_t^{obj}\| \|s_t^{obj}\|}$

126 Compared to the original loss function of FOCUS (defined in Appendix E), the world model loss
 127 becomes: $\mathcal{L}_{ocwm} = \mathcal{L}_{FOCUS} + \mathcal{L}_{pos}$

128 **Latent-Conditioned Policy Learning.** The introduction of the *latent positional encoder* enables
 129 the conditioning over the target object’s latent. By encoding a desired target position p_g^{obj} , the target
 130 object’s latent state s_g^{obj} is inferred. The latter serves as the conditioning factor for the policy network:
 131 $\pi_{LCP}(a_t|s_t, s_g^{obj})$. To incentivize the policy to move the entity of interest to the target location, we
 132 maximize the negative latent distance between \hat{s}_t^{obj} and s_g^{obj} . The distance function used is cosine
 133 similarity. r_{LCP} becomes then: $r_{LCP} = \frac{\hat{s}_t^{obj} \cdot s_g^{obj}}{\|\hat{s}_t^{obj}\| \|s_g^{obj}\|}$

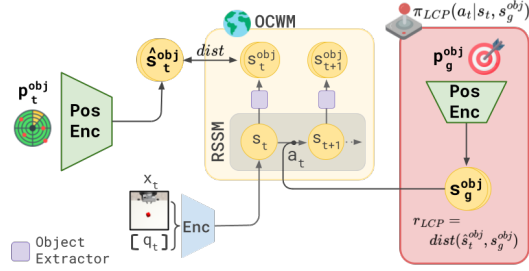


Figure 3: LCP leverages an object-centric representation. With the latent position encoder network, the agent learns to predict the latent of each object in the scene given the sole object position. The policy is then conditioned on an object latent target obtained from the target goal observation. Distance functions are expressed as cosine similarities.

134 **Visual targets.** Additionally with respect to PCP, LCP enables conditioning the policy on visual
 135 targets. In this case, the agent does not use the latent position encoder. Instead, given a visual
 136 observation representing the goal target position for the object, the world model can infer the
 137 corresponding world model state, using the encoder and the posterior. Then given such a state, the
 138 object extractor allows extracting the target latent state s_g^{obj} , which is used in the reward computation.
 139

140 4 Results

141 We now present the evaluation of the
 142 trained models (training details in Ap-
 143 pendix G) for a set of 4 environments
 144 (Appendix F). The score function considered is presented in Appendix, Eq. 2 .

	Dreamer	FOCUS	Dreamer w/ PCP	LEXA (cosine)	LEXA (temporal)	FOCUS w/ LCP
Reacher	0.26 ± 0.19	0.29 ± 0.19	0.8 ± 0.14	0.92 ± 0.04	0.42 ± 0.2	0.91 ± 0.04
Cube move	0.35 ± 0.08	0.35 ± 0.15	0.34 ± 0.08	0.37 ± 0.03	0.39 ± 0.1	0.61 ± 0.12
Shelf place	0.4 ± 0.14	0.3 ± 0.18	0.58 ± 0.13	0.4 ± 0.11	0.4 ± 0.11	0.65 ± 0.14
Pick&Place	0.26 ± 0.21	0.21 ± 0.2	0.48 ± 0.26	0.34 ± 0.24	0.34 ± 0.25	0.45 ± 0.29

Table 1: Average score for 100 goal points equally distributed over the workspace. Performance is averaged over 3 seeds, ± indicates the std. error.

145 **Spatial-coordinates goal specification.** By providing the different agents with goals uniformly
 146 distributed in the workspace we extract the overall performance of each method. Results are presented
 147 in Table 1. Overall, the FOCUS agent equipped with PCP or LCP gives the best performance,
 148 followed by Dreamer + PCP. In the "Shelf place" environment, the latent representation of LCP
 149 represents best. Given that the camera is further away from the scene, we believe the agent is
 150 better able to deal with the inaccuracies that come from the inaccurate position readings (bigger
 151 segmentation mask → better granularity in position).

152 **Visual goal specification.** An emergence prop-
 153 erty of FOCUS + LCP is the possibility to define
 154 goals via different modalities. The policy π_{LCP}
 155 can be conditioned on the goal object latent \hat{s}_g^{obj}
 156 coming from the encoding of the visual goal x_g .
 157 We compare our method with visual goal condi-
 158 tioning against LEXA cosine and temporal.
 159 The goal locations are provided to the simulator
 160 which renders the corresponding goal observa-
 161 tions by "teleporting" the object to the correct
 162 location. The agent is then asked to matched
 163 the visual goal, after resetting the environment.
 164 Results are shown in Fig. 4, where the positional
 conditioning results are shown for reference.

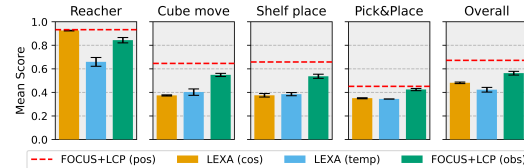


Figure 4: The mean score was achieved over 10 episodes with goal observations for latent conditioning. The performance of our method with spatial-coordinate goals (pos) is shown as a reference. Performance is averaged over 3 seeds.

165 LEXA matches the flat latent vector to the goal one. This proves helpful in the Reacher environment,
 166 where the only part that moves is the agent, and thus LEXA cosine achieves the best performance.
 167 LEXA cosine fails in the other tasks, given the presence of multiple entities in the observations and
 168 visual goals, i.e. the robotic arm and the object. where the model focuses on matching the visually
 169 predominant features i.e. the robotic arm. FOCUS+LCP performs better than both LEXA with cosine
 170 and temporal distance in all environments but the Reacher. When compared to the performance of
 171 FOCUS+LCP with spatial-coordinates goals, there is a decrease of ~10% in performance.

172 5 Conclusion

173 We analyzed the challenges in solving visual robotic positional tasks using generative world model-
 174 based agents. We found these systems suffer from information bottleneck issues when considering
 175 positional information for task resolution (i.e. goal position). The approaches we presented overcome
 176 this issue by providing the policy network with more direct access to the target information. Positional
 177 Conditioning Policy (PCP), allows direct conditioning on the target spatial coordinates. We showed
 178 PCP improves performance for any class of world models, including Dreamer-like "flat" world
 179 models and FOCUS-like object-centric world models. Latent Conditioning Policy (LCP), is an
 180 object-centric approach that we implement on top of FOCUS. This allows the conditioning of the
 181 policy on object-centric latent targets, enabling multimodal goal definition.

References

- 182
- 183 [1] S. Fujimoto, H. van Hoof, and D. Meger. Addressing function approximation error in actor-critic
184 methods, 2018.
- 185 [2] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy
186 deep reinforcement learning with a stochastic actor, 2018.
- 187 [3] D. Ha and J. Schmidhuber. World models. 2018. doi:10.5281/ZENODO.1207631. URL
188 <https://zenodo.org/record/1207631>.
- 189 [4] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi. Dream to control: Learning behaviors by latent
190 imagination. 2020. URL <https://arxiv.org/pdf/1912.01603.pdf>.
- 191 [5] S. Rajeswar, P. Mazzaglia, T. Verbelen, A. Piché, B. Dhoedt, A. Courville, and A. Lacoste.
192 Mastering the unsupervised reinforcement learning benchmark from pixels. 2023.
- 193 [6] D. Hafner, J. Pasukonis, J. Ba, and T. Lillicrap. Mastering diverse domains through world
194 models. *arXiv preprint arXiv:2301.04104*, 2023.
- 195 [7] P. Lancaster, N. Hansen, A. Rajeswaran, and V. Kumar. Modem-v2: Visuo-motor world models
196 for real-world robot manipulation, 2024.
- 197 [8] P. Wu, A. Escontrela, D. Hafner, K. Goldberg, and P. Abbeel. Daydreamer: World models for
198 physical robot learning, 2022.
- 199 [9] Y. Seo, J. Kim, S. James, K. Lee, J. Shin, and P. Abbeel. Multi-view masked world models for
200 visual robotic manipulation, 2023.
- 201 [10] S. Ferraro, P. Mazzaglia, T. Verbelen, and B. Dhoedt. Focus: Object-centric world models for
202 robotics manipulation, 2023.
- 203 [11] D. Hafner, T. P. Lillicrap, M. Norouzi, and J. Ba. Mastering atari with discrete world models.
204 In *ICLR*, 2021.
- 205 [12] N. Hansen, H. Su, and X. Wang. Td-mpc2: Scalable, robust world models for continuous
206 control, 2024.
- 207 [13] Y. Seo, D. Hafner, H. Liu, F. Liu, S. James, K. Lee, and P. Abbeel. Masked world models for
208 visual control, 2022.
- 209 [14] S. Tunyasuvunakool, A. Muldal, Y. Doron, S. Liu, S. Bohez, J. Merel, T. Erez, T. Lillicrap,
210 N. Heess, and Y. Tassa. dm_control: Software and tasks for continuous control. *Software*
211 *Impacts*, 6:100022, 2020. ISSN 2665-9638. doi:<https://doi.org/10.1016/j.simpa.2020.100022>.
212 URL <https://www.sciencedirect.com/science/article/pii/S2665963820300099>.
- 213 [15] Y. Zhu, J. Wong, A. Mandlekar, R. Martín-Martín, A. Joshi, S. Nasiriany, and Y. Zhu. robo-
214 suite: A modular simulation framework and benchmark for robot learning. In *arXiv preprint*
215 *arXiv:2009.12293*, 2020.
- 216 [16] H. Ma, J. Wu, N. Feng, C. Xiao, D. Li, J. Hao, J. Wang, and M. Long. Harmonydream: Task
217 harmonization inside world models, 2024. URL <https://arxiv.org/abs/2310.00344>.
- 218 [17] R. Mendonca, O. Rybkin, K. Daniilidis, D. Hafner, and D. Pathak. Discovering and achieving
219 goals via world models, 2021.
- 220 [18] D. P. Kingma and M. Welling. Auto-encoding variational bayes, 2022.
- 221 [19] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson. Learning latent
222 dynamics for planning from pixels. In *ICML*, pages 2555–2565, 2019.

- 223 [20] T. Yu, D. Quillen, Z. He, R. Julian, K. Hausman, C. Finn, and S. Levine. Meta-world: A
224 benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on*
225 *Robot Learning (CoRL)*, 2019. URL <https://arxiv.org/abs/1910.10897>.
- 226 [21] OpenAI, I. Akkaya, M. Andrychowicz, M. Chociej, M. Litwin, B. McGrew, A. Petron, A. Paino,
227 M. Plappert, G. Powell, R. Ribas, J. Schneider, N. A. Tezak, J. Tworek, P. Welinder, L. Weng,
228 Q. Yuan, W. Zaremba, and L. M. Zhang. Solving rubik’s cube with a robot hand. *ArXiv*,
229 [abs/1910.07113](https://arxiv.org/abs/1910.07113), 2019.
- 230 [22] K. Cobbe, C. Hesse, J. Hilton, and J. Schulman. Leveraging procedural generation to benchmark
231 reinforcement learning, 2020. URL <https://arxiv.org/abs/1912.01588>.
- 232 [23] J. Fan. A review for deep reinforcement learning in atari:benchmarks, challenges, and solutions,
233 2023. URL <https://arxiv.org/abs/2112.04145>.
- 234 [24] R. Sekar, O. Rybkin, K. Daniilidis, P. Abbeel, D. Hafner, and D. Pathak. Planning to explore
235 via self-supervised world models. In *ICML*, 2020.

236 Appendix

237 A Preliminaries

238 The agent is a robotic manipulator that, at each discrete timestep t receives an input x_t from the
239 environment. The goal of the agent is to move an object in the environment from its current position
240 p_t^{obj} to a target goal position p_g^{obj} .

241 In this work, we focus on observations composed of both visual and vector entities. Thus, $x_t = (o_t, v_t)$
242 is composed of the visual component o_t and of the vector v_t . The latter is a concatenation of
243 proprioceptive information of the robotic manipulator q_t , the object’s position p_t^{obj} , and the target
244 position p_g^{obj} . The target position can also be expressed through a visual observation x_g , from which
245 the agent should infer the corresponding p_g^{obj} to succeed in the positioning task.

246 A.1 Generative World Models

247 Generative world models learn a latent representation of the agent inputs using a variational auto-
248 encoding framework [18]. Dreamer-like agents [11, 6] implement the world model as a Recurrent
249 State-Space Model (RSSM) [19]. The encoder $f(\cdot)$ is instantiated as the concatenation of the
250 outputs of a CNN for high-dimensional observations and an MLP for low-dimensional proprioception.
251 Through the encoder network, the input x_t is mapped to an embedding e_t , which then is integrated
252 with dynamical information with respect to the previous RSSM state and the action taken a_t , resulting
253 in s_t features.

$$\begin{aligned} \text{Encoder: } & e_t = f(x_t) \\ \text{Posterior: } & p_\phi(s_{t+1}|s_t, a_t, e_{t+1}), \\ \text{Prior: } & p_\phi(s_{t+1}|s_t, a_t), \\ \text{Decoder: } & p_\theta(\hat{x}_t|s_t). \end{aligned}$$

254 Generally, the system either learns to predict the expected reward given the latent features [4], using a
255 reward predictor $p_\theta(\hat{r}_t|s_t)$. Alternatively, some world-model based methods adopt specialized ways
256 to compute rewards in imagination, as the goal-conditioned objectives in LEXA [17].

257 Rewards are computed on rollouts of latent states generated by the model and are used to learn the
258 policy π and value network v in imagination [4, 11, 6].

259 In our experiments, we consider a world model with a discrete latent space [11]. We also implement
260 advancements of the world model representation introduced in DreamerV3 [6], such as the application
261 of the symlog transform to the inputs, KL balancing, and free bits to improve the predictions of the
262 vector inputs and the robustness of the model.

263 A.2 Object-centric World Models

264 Compared to Dreamer-like *flat* world models, the world model of FOCUS [10] introduces the
265 following object-centric components:

$$\begin{aligned} \text{Object latent extractor: } & p_\theta(s_t^{obj}|s_t, c^{obj}), \\ \text{Object decoder: } & p_\theta(\hat{x}_t^{obj}, \hat{m}_t^{obj}|s_t^{obj}). \end{aligned}$$

266 Here, $x_t^{obj} = (o_t^{obj}, p_t^{obj})$ represents the object-centric inputs and it is composed of segmented RGB
267 images o_t^{obj} and object positions p_t^{obj} . The variable c^{obj} indicates which object is being considered.

268 Thanks to the *object latent extractor* unit, object-specific information is separated into distinct
269 latent representations s_t^{obj} . Two decoding units are present. The introduced object-centric decoder
270 $p_\theta(\hat{x}_t^{obj}, \hat{m}_t^{obj}|s_t^{obj})$ reconstructs each object’s related inputs x_t^{obj} and segmentation mask m_t^{obj} . The
271 original Dreamer-like decoder takes care of the reconstruction of the remaining vector inputs, i.e.
272 proprioception q_t and given goal targets p_g^{obj} .

273 We provide additional descriptions of the world model and policy learning losses, hyperparameters,
 274 and training details in the Appendix.

275 A.3 Object Positioning Tasks

276 In general terms, we consider positioning tasks the ones where an entity of interest has to be moved
 277 to a specific location. Two positioning scenarios are considered in this analysis: *pose reaching* and
 278 *object positioning*. Pose-reaching tasks can be seen as simplified positioning tasks where the entity
 279 of interest is part of the robotic manipulator itself. Pose-reaching tasks are interesting because these
 280 only require the agent to have knowledge of the proprioceptive information to infer their position in
 281 space and reach a given target. When interacting with objects instead, there is the additional necessity
 282 of knowing the position of the object entity in the environment. Then, the agent needs to be able to
 283 manipulate and move the entity to the provided target location.

284 For object positioning tasks, especially when considering a real-world setup, there is a significant
 285 advantage in relying mainly on visual inputs. It is convenient because it avoids the cost and difficulty
 286 associated with tracking additional state features, such as the geometrical shape of objects in the
 287 scene or the presence of obstacles. Some synthetic benchmarks additionally make use of "virtual"
 288 visual targets for positioning tasks [14, 20], which strongly facilitates the learning of these tasks,
 289 leveraging rendering in simulation. However, applying such "virtual" targets in real-world settings is
 290 not often feasible. Non-visual target locations can be provided as spatial coordinates. Alternatively,
 291 an image showing the target location could be used to specify the target’s position.

292 **Rewards and evaluation criteria.** When applying RL algorithms to a problem, a heavily engineered
 293 reward function is generally necessary to guide the agent’s learning toward the solution of the task
 294 [21]. The object positioning setup allows us to consider a natural and intuitive reward definition
 295 that scales across different agents and environments. We define the reward as the negative distance
 296 between the position of the entity of interest and the goal target position:

$$r_t = -\text{distance}(\text{object}, \text{target}) = -\|p_t^{obj} - p_g^{obj}\|_2. \quad (1)$$

297 In the spirit of maintaining a setup that is as close as possible to a real-world one, to retrieve positional
 298 information p_t of the objects we rely on image segmentation information, rather than using the
 299 readings provided from the simulator. For each entity of interest, the related position is extracted by
 300 computing the centroid of the segmentation mask and subsequently transformed according to the
 301 camera extrinsic and intrinsic matrices to obtain the absolute position with respect to the workspace.

302 For evaluation purposes, we use the goal-normalized score function:

$$\text{normalized score} = \exp\left(-\frac{\|p_t^{obj} - p_g^{obj}\|_2}{\|p_g^{obj}\|_2}\right) \quad (2)$$

303 As detailed in the Appendix, the above function allows us to rescale performance between 0 and 1,
 304 where 1 = expert performance, a common evaluation strategy in RL [22, 23].

305 B Normalized score

306 Scaling performance using expert performance is a common evaluation strategy in RL [22, 23]. In
 307 our problem, we define the reward as the negative distance:

$$r_t = -r(p_t^{obj}) = -\|p_t^{obj} - p_g^{obj}\|_2. \quad (3)$$

308 For a given goal p_g^{obj} , $r_t \in]-\infty, 0]$. In order to compare different tasks, where distances may
 309 have different magnitudes, we divide the rewards r_t by the typical reward range. This is given by
 310 $r_{max} - r_{min}$, where $r_{min} = r(p_0^{obj})$, with p_0 being the initial position of the object (this is normally
 311 around the origin, and $r_{max} = r(p_g^{obj}) = 0$).

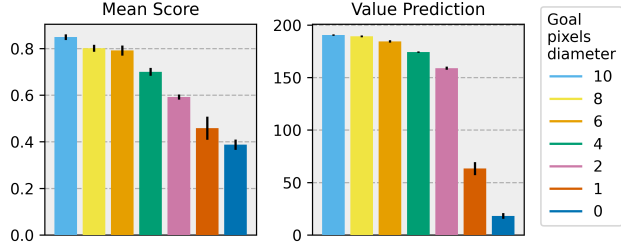


Figure 5: Dreamer virtual visual goal modulation experiments on the Reacher environment. Value prediction from the value network is shown to highlight the policy’s awareness of the lack of information with respect to the target goal.

312 Thus, we obtain:

$$s_t = r_t / (r_{max} - r_{min}) \quad (4)$$

$$= r(p_t^{obj}) / (0 - r(p_0^{obj})) = \quad (5)$$

$$= -\|p_t^{obj} - p_g^{obj}\|_2 / (0 + \|0 - p_g^{obj}\|_2) \quad (6)$$

$$= -\|p_t^{obj} - p_g^{obj}\|_2 / \|p_g^{obj}\|_2 \quad (7)$$

313 Finally, we apply the exp operator, to make values positive and bring them in the $[0, 1]$ range, where
 314 1 is the expert score:

$$\text{normalized score} = \exp\left(-\frac{\|p_t^{obj} - p_g^{obj}\|_2}{\|p_g^{obj}\|_2}\right) \quad (8)$$

315 C Target size ablation

316 In Figure 5, we present a study where the Dreamer model is trained on the Reacher environment with
 317 varying visual target sizes.

318 We observe that the reduction in pixel information regarding the target adversely affects the target
 319 representation within the model, resulting in a deficiency of this information being conveyed to the
 320 policy network. The policy struggles to learn to position the entity at the correct location, and we
 321 observe that this is correctly reflected in the value function’s predictions. This means the policy
 322 is aware that is not being able to reach the goal. With small targets (< 5 pixels diameters), the
 323 representation tends to put more attention on other visually predominant aspects of the environment,
 324 struggling to predict the position of the target. In the case of a single pixel target, the amount of target
 325 information equals the one of a positional vector and, as expected, the task performance is equally
 326 low.

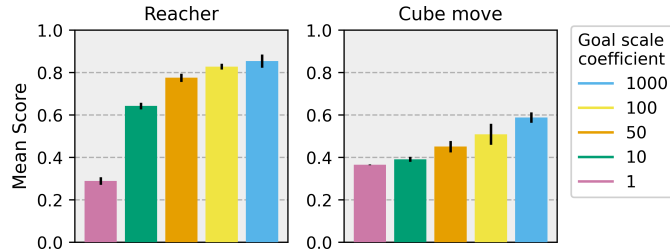


Figure 6: Dreamer trained with goal scaling modulation on the Reacher and Cube move environments.

327 **D Loss rescaling ablation**

328 To overcome the identified information bottleneck, different strategies can be considered. The
 329 simplest one is the re-scaling of the loss components in the decoder to incentivize the model’s
 330 encoding of the target information. This approach requires finding the optimal scaling factor between
 331 the different decoding components, given the complexity of the environment at hand (i.e. 2D or 3D)
 332 and the amount of relevant pixels. In Figure 6, we present supporting experiments based on Dreamer,
 333 where we vary the importance of the target in the loss of the world model, using different coefficients.
 334 We observe that very high coefficients improve the target’s reconstruction and thus allow the agent to
 335 learn the task. However, the optimal loss coefficient may vary, depending on the complexity of the
 336 environment and the presence of information-rich observations. As this naive solution may require
 337 extensive hyperparameter tuning for each new scenario, we aim to find more robust strategies for
 338 overcoming this issue.

339 **E FOCUS objective**

340 Training of the FOCUS architecture is guided by the following loss function:

$$\mathcal{L}_{\text{FOCUS}} = \mathcal{L}_{\text{dyn}} + \mathcal{L}_{\text{state}} + \mathcal{L}_{\text{obj}}. \tag{9}$$

341 \mathcal{L}_{dyn} refers to the dynamic component of the RSSM, and equals too:

$$\mathcal{L}_{\text{dyn}} = D_{\text{KL}}[p_{\phi}(s_{t+1}|s_t, a_t, e_{t+1})||p_{\phi}(s_{t+1}|s_t, a_t)]. \tag{10}$$

342 the backpropagation is balanced and clipped below 1 nat as in DreamerV3 [6].

343 The object loss component is instantiated as the composition of NLL over the mask and RGB mask
 344 reconstructions:

$$\mathcal{L}_{\text{obj}} = -\log \underbrace{p(\hat{m}_t)}_{\text{mask}} - \log \sum_{\text{obj}=0}^N \underbrace{m_t^{\text{obj}} p_{\theta}(\hat{x}_t^{\text{obj}}|s_t^{\text{obj}})}_{\text{masked reconstruction}} \tag{11}$$

345 Finally, the decoder learns to reconstruct the rest of vector state information v_t by minimization of
 346 the negative log-likelihood (NLL) loss:

$$\mathcal{L}_{\text{state}} = -\log p_{\theta}(\hat{q}_t, p_g^{\text{obj}}|s_t) \tag{12}$$

347 **F Baselines and Environments**

348 For the evaluation of the proposed method we consider several manipulation environments (Figure 7):

- 349 • **Reacher** (DMControl): which, as described previously, represents a pose-reaching position-
 350 ing task.
- 351 • **Cube move** (Robosuite): where considered target locations are on the 2D plane of the table,
 352 no height placement is considered.

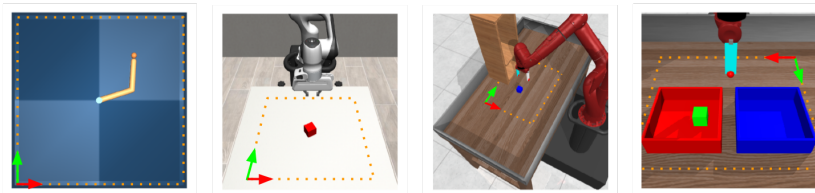


Figure 7: Simulation environments with relative workspace, delimited by an orange dotted line, and the reference frames indicated with arrows.

353 • **Shelf place** and **Pick&Place** (Metaworld): The robotic manipulator has to place the cube
354 at the given target location. Considered target locations are on the 2D space in front of the
355 robotic arm.

356 In all environments, the reward signal is defined as the distance between the entity of interest (in the
357 Reacher environment, this is the end-effector) and the target location. All considered environments
358 lack any visual target; the target is provided as an input vector containing spatial coordinates.

359 We benchmark our methods against various baselines:

360 • **Dreamer**: based on a PyTorch DreamerV2 implementation, but integrated with input vector
361 symlog transformation and KL balancing of the latent dynamic representation, from the
362 DreamerV3 paper.

363 • **FOCUS**: An object-centric world model implementation based on DreamerV2, also inte-
364 grated with input vector symlog transformation and KL balancing of the latent dynamic
365 representation.

366 • **LEXA**: Based on DreamerV2, this is a latent goal-conditioned method. The conditioning is
367 based on the full latent target. Both proposed distance methods (cosine and temporal) are
368 considered. We adopted our own PyTorch implementation for LEXA.

369 G Training details and Hyperparameters

370 All methods are trained following an offline RL training scheme. The offline datasets contain 1M
371 steps in the environment, which are collected using the object-centric exploration strategy proposed
372 in [10]. The datasets are loaded in the replay buffer of the offline agents, and the training is conducted
373 for 250K steps. Both world model and agent are updated at every training step. V100-16GB GPUs
374 have been used for all experiments. Our proposed methods (i.e. Dreamer/FOCUS + PCP, FOCUS +
375 LCP) took roughly 18 hours to complete each training run.

376 The hyperparameters used for the main implementation of the world models and agent are the same
377 used in DreamerV2 [11] official implementation. Symlog function is applied at every input. KL
378 balancing as in DreamerV3 [6] is implemented.

379 With reference to FOCUS model, we have the following additional parameters:

380 • Object-extractor: MLP composed of 2 layers, 512 units, ReLU activation;

381 With reference to FOCUS + LCP model, we have the following additional parameters:

382 • Object-encoder: MLP composed of 4 layers, 400 units, ReLU activation;

383 • Distance method object-encoder objective: Cosine similarity (also tested MSE)

384 • Distance method actor policy objective: Cosine similarity (also tested MSE)

385 H Heatmaps positioning tasks

386 To highlight the performance distribution over the different goals in the environment, in Fig. 8 we
387 present heatmaps with the score function for each target location in the workspace. Results are
388 presented for all the different tasks. As expected, both Dreamer and FOCUS have poor performances,
389 resulting in only a few positions being reached with a high score. All the proposed methods have a
390 similar distribution, reaching goals spread all over the environment.

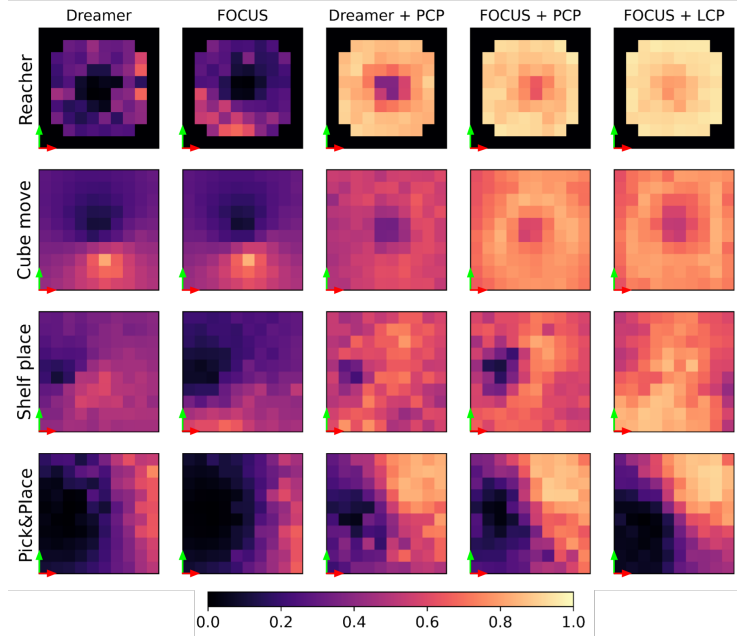


Figure 8: Heatmaps of the mean achieved score for uniformly spread targets in the workspace. References frames refers to the one presented in the figures of Table 1. The score notation is expressed as the notation presented in Eq. 2. Results are averaged over 3 seeds.

391 I Offline Training Curves

392 Offline training curves are presented in Figure 9. In general FOCUS + PCP/LCP have faster
 393 convergence when compared to all other methods. Only for the Reacher environment, LEXA cosine
 394 converge faster.

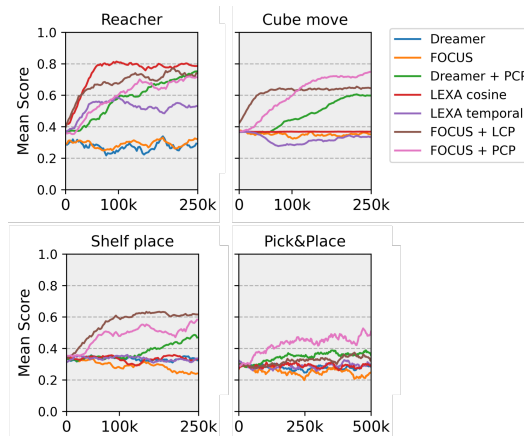


Figure 9: Offline training curves. Standard deviation is omitted for graphical reasons. Mean score refers to eq. 2 and is computed over 5 evaluation episodes, performed during the offline training. For each episode, a random goal is selected out of a pool of 10 manually engineered ones.

395 J Explorations strategies

396 In the presented work each model is trained offline from a pre-recorded dataset. The dataset of
 397 choice is obtained from pure exploration behavior. In Fig. 10 we compare the general performance

398 of LCP when trained on datasets acquired using different exploration strategies. We consider the
399 object-centric entropy maximization method proposed by Ferraro et al. [10] and Plan2Explore [24].

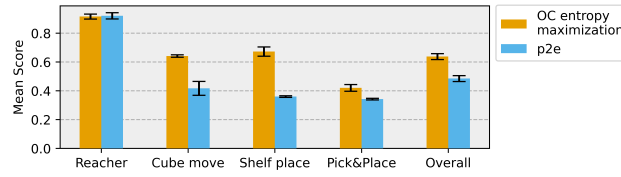


Figure 10: Mean score achieved over 10 episodes for models trained with both datasets obtained from FOCUS exploration method (Object-Centric entropy maximization) and Plan2Explore. The score is expressed according to equation 2.

400 Overall exploring by maximizing the entropy over the object’s latent, gives better performance in
401 the downstream task. We hypothesize this is related to the focus the exploration strategy puts on the
402 object of interest while disregarding background aspects in the scene.