

Object-based active inference

Ruben S. van Bergen and Pablo L. Lanillos

Department of Artificial Intelligence &
Donders Institute for Brain, Cognition & Behavior
Radboud University, Nijmegen, the Netherlands
{ruben.vanbergen, pablo.lanillos}@donders.ru.nl

Abstract The world consists of objects: distinct entities possessing independent properties and dynamics. For agents to interact with the world intelligently, they must translate sensory inputs into the bound-together features that describe each object. These object-based representations form a natural basis for planning behavior. Active inference (AIF) is an influential unifying account of perception and action, but existing AIF models have not leveraged this important inductive bias. To remedy this, we introduce ‘object-based active inference’ (OBAI), marrying AIF with recent deep object-based neural networks. OBAI represents distinct objects with separate variational beliefs, and uses selective attention to route inputs to their corresponding object slots. Object representations are endowed with independent action-based dynamics. The dynamics and generative model are learned from experience with a simple environment (active multi-dSprites). We show that OBAI learns to correctly segment the action-perturbed objects from video input, and to manipulate these objects towards arbitrary goals.

Keywords: Multi-object representation learning · Active inference

1 Introduction

Intelligent agents are not passive entities that observe the world and learn its causality. They learn the relationship of action and effect by interacting with the world, in order to fulfil their goals [1]. In higher-order intelligence, such as exhibited by primates, these interactions very often take place at the level of objects [2,3]. Whether picking a ripe fruit from a tree branch, kicking a football, or taking a drink from a glass of water; all require reasoning and planning in terms of objects. Objects, thus, are natural building blocks for representing the world and planning interactions with it.

While there have been recent advances in unsupervised multi-object representation learning and inference [4,5], to the best of the authors knowledge, no existing work has addressed how to leverage the resulting representations for generating actions. In addition, object perception itself could benefit from being placed in an active loop, as carefully selected actions could resolve ambiguity about object properties (including their segmentations - i.e., which inputs belong to which objects). Meanwhile, state-of-the-art behavior-based learning

(control), such as model-free reinforcement learning [6] uses complex encoding of high-dimensional pixel inputs without taking advantage of objects as an inductive bias (though see [7,8]).

To bridge the gap between these different lines of work, we here introduce ‘object-based active inference’ (OBAI, pronounced /əˈbeɪ/), a new framework that combines deep, object-based neural networks [4] and active inference [9,10]. Our proposed neural architecture functions like a Bayesian filter that iteratively refines perceptual representations. Through selective attention, sensory inputs are routed to high-level object modules (or *slots* [5]) that encode each object as a separated probability distribution, whose evolution over time is constrained by an internal model of action-dependent object dynamics. These object representations are highly compact and abstract, thus enabling efficient unrolling of possible futures in order to select optimal actions in a tractable manner. Furthermore, we introduce a closed-form procedure to learn preferences or goals in the network’s latent space.

As a proof-of-concept, we evaluate our proposed framework on an active version of the multi-dSprites dataset, developed for this work (See Fig. 1a). Our preliminary results show that OBAI is able to: *i*) learn to segment and represent objects *ii*) learn the action-dependent, object-based dynamics of the environment; and *iii*) plan in the latent space – obviating the need to imagine detailed pixel-level outcomes in order to generate behavior. This work is a first step towards building more complex object-based active inference systems that can perform more cognitively challenging tasks on naturalistic input.

2 Methods

2.1 Object-structured generative model

We extend the IODINE architecture proposed in [4] for object representation learning, to incorporate dynamics. Zablotskaia et al. [11] previously developed a similar extension to IODINE, in which object dynamics were modeled implicitly, through LSTM units operating one level below the latent-space representation. Here, we instead implement the dynamics directly in the latent space, and allow these dynamics to be influenced by actions on the part of the agent.

Like IODINE, our framework relies on iterative amortized inference [12] (IAI) on an object-structured generative model. This model describes images of up to K objects with a Normal mixture density (illustrated in Fig. 1):

$$p(o_i|\{\mathbf{s}^{(k)}\}_{k \in 1:K}, m_i) = \sum_k [m_i = k] \mathcal{N}(g_i(\mathbf{s}^{(k)}), \sigma_o^2) \quad (1)$$

where o_i is the value of the i -th image pixel, $\mathbf{s}^{(k)}$ is the state of the k -th object, $g_i(\bullet)$ is a decoder function (implemented as a deep neural network (DNN)) that translates an object state to a predicted mean value at pixel i , σ_o^2 is the variability of pixels around their mean values and, crucially, m_i is a categorical variable that

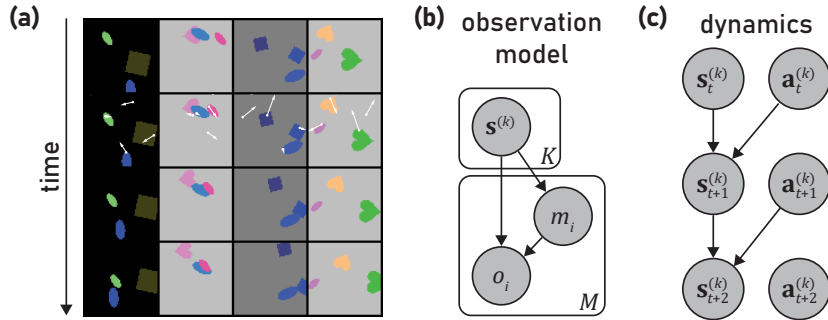


Figure 1. Environment and generative model. (a) Active multi-dSprites. Non-zero accelerations in action fields (in the 2nd frame) are indicated as white arrows originating at the accelerated grid location. (b) Object-structured generative model for a single image (time indices omitted here for clarity of exposition). (c) Dynamics model in state-space for a single object k , shown for three time points.

indicates which object (out of a possible K choices) pixel i belongs to¹. Note that the same decoder function is shared between objects. The pixel assignments themselves also depend on the object states:

$$p(m_i|\{\mathbf{s}^{(k)}\}_{k \in 1:K}) = \text{Cat}\left(\text{Softmax}\left(\{\pi_i(\mathbf{s}^{(k)})\}_{k \in 1:K}\right)\right) \quad (2)$$

where $\pi_i(\bullet)$ is another DNN that maps an object state to a log-probability at pixel i , which (up to a constant of addition) defines the probability that the pixel belongs to that object. Marginalized over the assignment probabilities, the pixel likelihoods are given by:

$$p(o_i|\{\mathbf{s}^{(k)}\}_{k \in 1:K}) = \sum_k \hat{m}_{ik} \mathcal{N}\left(g_i(\mathbf{s}^{(k)}), \sigma_o^2\right) \quad (3)$$

$$\hat{m}_{ik} = p(m_i = k|\{\mathbf{s}^{(k)}\}_{k \in 1:K}) \quad (4)$$

During inference, the soft pixel assignments $\{\hat{m}_{ik}\}$ introduce dynamics akin to selective attention, as each object slot is increasingly able to focus on those pixels that are relevant to that object.

2.2 Incorporating action-dependent dynamics

So far, this formulation is identical to the generative model in IODINE. We now extend this with an action-based dynamics model. We want to endow objects with (approximately) linear dynamics, and to allow actions that accelerate the objects. First, we redefine the state of an object at time point t in generalized coordinates, i.e. $\mathbf{s}_t^\dagger = \begin{bmatrix} \mathbf{s}_t \\ \mathbf{s}'_t \end{bmatrix}$, where \mathbf{s}' refers to the first-order derivative of the

¹ Note the use of Iverson-bracket notation; the bracket term is binary and evaluates to 1 iff the expression inside the brackets is true.

state. The action-dependent state dynamics are then given by:

$$\mathbf{s}_t^{(k)} = \mathbf{s}_{t-1}^{(k)} + \mathbf{D}\mathbf{a}_{t-1}^{(k)} + \sigma_s \epsilon_1 \quad (5)$$

$$\mathbf{s}_t^{(k)} = \mathbf{s}_{t-1}^{(k)} + \mathbf{s}_t^{\prime(k)} + \sigma_s \epsilon_2 \quad (6)$$

where $\mathbf{a}_t^{(k)}$ is the action on object k at time t . This action is a 2-D vector that specifies the acceleration on the object in pixel coordinates. Multiplication by \mathbf{D} (which is learned during training) transforms the pixel-space acceleration to its effect in the latent space². Equations 5-6 thus define the object dynamics model $p(\mathbf{s}_t^{\dagger(k)} | \mathbf{s}_{t-1}^{\dagger(k)}, \mathbf{a}_{t-1}^{(k)})$.

We established that $\mathbf{a}_t^{(k)}$ is the action on object k in the model at a given time. However, note that the correspondence between objects represented by the model, and the true objects in the (simulated) environment, is unknown.³ To solve this correspondence problem, we introduce the idea of *action fields*. An action field $\Psi = [\psi_1, \dots, \psi_M]^T$ is an $[M \times 2]$ matrix (with M the number of pixels in an image or video frame), such that the i -th row in this matrix (ψ_i) specifies the (x,y)-acceleration applied at pixel i . In principle, a different acceleration can be applied at each pixel coordinate (in practice, we apply accelerations sparsely). These pixel-wise accelerations affect objects through the rule that each object receives the sum of all accelerations that occur at its visible pixels:

$$\mathbf{a}_t^{(k)} = \sum_i [m_i = k] \psi_i + \sigma_\psi \epsilon_3 \quad (7)$$

where we include a small amount of Normally distributed noise in order to make this relationship amenable to variational inference. This definition of actions in pixel-space is unambiguous and allows the model to interact with the environment.

2.3 Inference

On the generative model laid out in the previous section, we perform *iterative amortized inference* (IAI). IAI generalizes variational autoencoders (VAEs), which perform inference in a single feedforward pass, to architectures which use several iterations (implemented in a recurrent network) to minimize the Evidence Lower Bound (ELBO). As in VAEs, the final result is a set of variational beliefs in the latent space of the network. In our case, this amounts to inferring

² Since the network will be trained unsupervised, we do not know in advance the nature of the latent space representation that will emerge. In particular, we do not know in what format (or even if) the network will come to represent the positions of the objects.

³ In particular, since the representation across objects slots in the network is permutation-invariant, their order is arbitrary – just as the order in the memory arrays that specify the environment is also arbitrary. Thus, the object-actions, as represented in the model, cannot be unambiguously mapped to objects in the environment that the agent interacts with. This problem is exacerbated if the network has not inferred the object properties and segmentations with perfect accuracy or certainty, and thus cannot accurately or unambiguously refer to a true object in the environment.

$q(\{\mathbf{s}^{\dagger(k)}, \mathbf{a}^{(k)}\}_{k \in 1:K})$. We choose these beliefs to be independent Normal distributions. Inference and learning both minimize the following ELBO loss:

$$\begin{aligned} \mathcal{L} = & - \sum_{t=0}^T \left[\mathcal{H} \left(q \left(\{\mathbf{s}_t^{\dagger(k)}, \mathbf{a}_t^{(k)}\} \right) \right) + E_{q(\{\mathbf{s}_t^{(k)}\})} [\log p(\mathbf{o}_t | \{\mathbf{s}_t^{(k)}\})] \right. \\ & \left. + \sum_k E_{q(\mathbf{a}_t^{(k)})} [\log p(\mathbf{a}_t^{(k)} | \Psi_t)] + \sum_k E_{q(\mathbf{s}_t^{\dagger(k)}, \mathbf{s}_{t-1}^{\dagger(k)}, \mathbf{a}_{t-1}^{(k)})} [\log p(\mathbf{s}_t^{\dagger(k)} | \mathbf{s}_{t-1}^{\dagger(k)}, \mathbf{a}_{t-1}^{(k)})] \right] \quad (8) \end{aligned}$$

for some time horizon T . Note that for $t = 0$, we define $p(\mathbf{s}_t^{\dagger(k)} | \mathbf{s}_{t-1}^{\dagger(k)}, \mathbf{a}_{t-1}^{(k)}) = p(\mathbf{s}^{\dagger(k)}) = \prod_{jk} \mathcal{N}(s_j^{\dagger(k)}; 0, 1)$, i.e. a fixed standard-Normal prior. To compute $E_{q(\mathbf{a}_t^{(k)})} [\log p(\mathbf{a}_t^{(k)} | \Psi_t)]$, we employ a sampling procedure, described in Appendix D.

The IAI architecture consists of a decoder module that implements the generative model, and a refinement module which outputs updates to the parameters λ of the variational beliefs. Mirroring the decoder, the refinement module consists of K copies of the same network (sharing the same parameters), such that refinement network k outputs updates to $\lambda^{(k)}$. Network architectures for the decoder and refinement modules are detailed in Appendix B. To perform inference across multiple video frames, we simply copy the refinement and decoder networks across frames as well as object slots. Importantly, as in [4], each refinement network instance also receives as input a stochastic estimate of the current gradient $\nabla_{\lambda_t^{(k)}} \mathcal{L}$. Since the ELBO loss includes a temporal dependence term between time points, the inference dynamics in the network are automatically coupled between video frames, constraining the variational beliefs to be consistent with the dynamics model. To infer $q(\{\mathbf{a}^{(k)}\}_{k \in 1:K})$, we employ a separate (small) refinement network (again copied across objects slots; details in Appendix B.2).

2.4 Task & training

We apply OBAI to a simple synthetic environment, developed for this work, which we term *active-dSprites*. This environment was created by re-engineering the original dSprites shapes [13], to allow these objects to be translated at will and by non-integer pixel offsets. The environment simulates these shapes moving along linear trajectories that can be perturbed through the action fields we introduced above. In the current work, OBAI was trained on pre-generated experience with this environment, with action fields sampled arbitrarily (i.e. not based on any intent on the part of the agent).

Specifically, we generated video sequences (4 frames, 64×64 pixels) from the active-dSprites environment, with 3 objects per video, in which we applied an action field only at a single time point (in the 2nd frame). Action fields were sparsely sampled such that (whenever possible) every object received exactly one non-zero acceleration at one of its pixels. Exactly one background pixel also received a non-zero acceleration, to encourage the model to learn not to assign background pixels to the segmentation masks of foreground objects. In practice,

the appearance of the background was unaffected by these actions (conceptually, the background can be thought of as an infinitely large plane extending outside the image frame, and thus shifting it by any amount will not change its visual appearance in the image).

OBAI was trained on 50,000 pre-generated video sequences, to minimize the ELBO loss from equation 8. This loss was augmented to include the losses at intermediate iterations of the inference procedure, and this composite loss was backpropagated through time to compute parameter gradients. More details about the environment and training procedure can be found in Appendices A & C.

2.4.1 Learning goals in the latent space

The active-dSprites environment was conceived to support cognitive tasks that require object-based reasoning. A natural task objective in active-dSprites is to move a certain object to a designated location. This type of objective is simple in and of itself, but the rules that determine which object must be moved where can be arbitrarily complex. For now, we restrict ourselves to the simple objective of moving all objects in a scene to a fixed location, and focus on how to encode this objective. We follow previous Active Inference work (e.g. [14,15]) in conceptualizing goals as a preference distribution \tilde{p} . However, rather than defining this preference to be over observations, as is common (though see [16,17,18]), we instead opt to define it over latent states, i.e. $\tilde{p}(\{\mathbf{s}^\dagger(k)\})$, which simplifies action selection (a full discussion of the merits of this choice is outside the scope of this paper).

Assuming that we can define a preference over the true state of the environment, \mathbf{s}_{true} (e.g. the ground-truth object positions), the preference distribution in latent space can be obtained through the following importance-sampling procedure:

$$\tilde{p}(\mathbf{s}) \propto \sum_j p(\mathbf{s}|\mathbf{o}_j^*)u_j \approx \sum_j q(\mathbf{s}|\mathbf{o}_j^*)u_j \quad (9)$$

$$\mathbf{o}_j^* \sim p(\mathbf{o}|\mathbf{s}_{\text{true}_j}^*), \quad u_j = \tilde{p}(\mathbf{s}_{\text{true}_j}^*), \quad \mathbf{s}_{\text{true}_j}^* \sim p(\mathbf{s}_{\text{true}}) \propto \text{Constant} \quad (10)$$

This allows the latent-space preference to be estimated in closed form from a set of training examples, constructed by sampling true states uniformly from the environment and rendering videos from these states. Inference is performed on the resulting videos, and the latent-state preference is computed as the importance-weighted average of the inferred state-beliefs (alternatively, we can sample states directly from $\tilde{p}(\mathbf{s}_{\text{true}}^*)$, and let the importance weights drop out of the equation). In particular, if $\tilde{p}(\mathbf{s}_{\text{true}}^*)$ is Normal, then the preference in the latent space is also Normal:

$$q(\mathbf{s}|\mathbf{o}_j^*) = \mathcal{N}(\boldsymbol{\mu}(\mathbf{o}_j^*), \boldsymbol{\sigma}(\mathbf{o}_j^{*2})), \quad \tilde{p}(\mathbf{s}) = \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\sigma}}^2) \quad (11)$$

$$\tilde{\boldsymbol{\mu}} = \frac{1}{\sum_j u_j} \sum_j u_j \boldsymbol{\mu}(\mathbf{o}_j^*), \quad \tilde{\boldsymbol{\sigma}} = \sqrt{\frac{1}{\sum_j u_j} \sum_j u_j ((\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}(\mathbf{o}_j^*))^2 + \boldsymbol{\sigma}(\mathbf{o}_j^{*2}))} \quad (12)$$

2.4.2 Planning actions

OBAI can plan actions aimed at bringing the environment more closely in line with its learned preferences. Specifically, we choose actions that minimize the Free Energy of the Expected Future (FEEF) [18]. When preferences are defined with respect to latent states, the FEEF of a policy (action sequence) $\pi = \left\{ [\mathbf{a}_1^{(k)}, \mathbf{a}_2^{(k)}, \dots, \mathbf{a}_T^{(k)}] \right\}_{k \in 1:K}$ is given by:

$$\mathcal{G}(\pi) = \sum_{\tau=1}^T \sum_k D_{KL} \left(q(\mathbf{s}_\tau^{(k)} | \pi) \| \tilde{p}(\mathbf{s}) \right) \quad (13)$$

where $q(\mathbf{s}_\tau^{(k)} | \pi)$ is the policy-conditioned variational prior, obtained by propagating the most recent state beliefs through the dynamics model by the requisite number of time steps.

Given this objective, the optimal policy can be calculated in closed form for arbitrary planning horizons. In this work, as a first proof-of-principle, we only consider greedy, one-step-ahead planning. In this case, the optimal "policy" (single action per object) is given by:

$$\hat{\mathbf{a}}^{(k)} = (\mathbf{D}^T \mathbf{L} \mathbf{D})^{-1} \mathbf{D}^T \mathbf{L} (\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}_s^{(k)}) \quad (14)$$

where $\mathbf{L} = \text{diag}(\tilde{\boldsymbol{\sigma}}^{-2})$, and $\boldsymbol{\mu}_s^{(k)}$ is the mean of the current state belief for object k .

3 Results

3.1 Object segmentation and reconstruction in dynamic scenes

We first evaluated OBAI on its inference and reconstruction capabilities, when presented with novel videos of moving objects (not seen during training). To evaluate this, we examined the quality of its object segmentations, and of the video frame reconstructions (Fig. 2). Segmentation quality was computed using the Adjusted Rand Index (ARI), as well as a modified version of this index that only considers (ground-truth) foreground pixels (FARI). Across a test set of 10,000 4-frame video sequences of 3 randomly sampled moving objects each, OBAI achieved an average ARI and FARI of 0.948 and 0.939, respectively (where 1 means perfect accuracy and 0 equals chance-level performance), and a MSE of 9.51×10^{-4} (note that pixel values were in the range of $[0, 1]$). For comparison, a re-implementation of IODINE, trained on 50,000 static images of 3 dSprite objects, achieved an ARI of 0.081, FARI of 0.856 and MSE of 1.63×10^{-3} (on a test set of identically sampled static images). The very low ARI score reflects the fact that IODINE has no built-in incentive to assign background pixels to their own object slot. OBAI, on the other hand, has to account for the effects of actions being applied to the background, which must not affect the dynamics of the foreground objects. Thus, for OBAI to accurately model the dynamics in the training data, it must learn not to assign background pixels to the segmentation masks of foreground objects, lest an action might be placed on one of these spurious pixels.

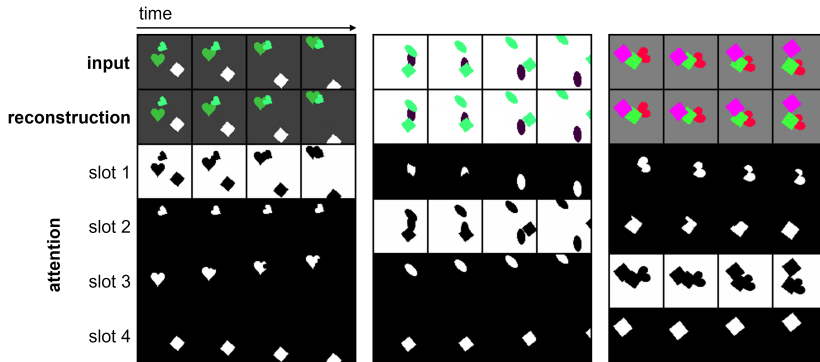


Figure 2. Reconstruction and segmentation of videos of moving objects with actions. Three instances of segmentation from the test set of videos. The masks show how each slot attends to the three objects and the background. The action field in this experiment is randomly generated for each instance at the beginning of the simulation.

	ARI (\uparrow)	F-ARI (\uparrow)	MSE (\downarrow)
IODINE (static)	0.081	0.856	1.63×10^{-3}
OBAI (ours; 4 frames)	0.948	0.939	9.51×10^{-4}

Table 1. Quantitative segmentation and reconstruction results.

3.2 Predicting the future state of objects

An advantage of our approach is that the network can predict future states of the world at the level of objects using the learned state dynamics. Figure 3 shows three examples of the network predicting future video frames. The first 4 video frames are used by the network to infer the state. Afterwards, we extrapolate the inferred state and dynamics of the last observed video frame into the future, and decode the thus-predicted latent states into predicted video frames. These predictions are highly accurate, with a MSE of 4×10^{-3} .

3.3 Goal-directed action planning

Can OBAI learn and accomplish behavioral objectives? As a first foray into this question, we asked OBAI to learn fixed preference distributions defined in the true state-space of the environment, using the method described in section 2.4.1. Specifically, we placed a Gaussian preference distribution on the location of the object and had the network learn the corresponding preference in its latent space from a set of 10,000 videos (annotated with the requisite importance weights). We then presented the network with static images of dSprite objects in random locations, and asked it to "imagine" the action that would bring the state of the environment into alignment with the learned preference. Finally, we applied this

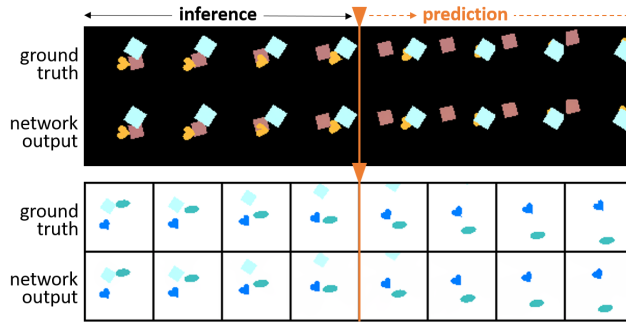


Figure 3. Prediction of future video frames. Two instances of prediction from the test set of videos. We let the network perform inference on 4 consecutive frames, and then predict the future.

imagined action to the latent state, and decoded the image that would result from this. As illustrated in Figure 4, the network is reliably able to imagine actions that would accomplish the goals we wanted it to learn.

4 Conclusion

This work seeks to bridge an important gap in the field. On the one hand, computer vision research has developed object-based models, but these only perform passive inference. On the other hand, there is a wealth of research on behavioral learning (e.g. reinforcement learning and active inference), which has not generally leveraged objects as an inductive bias, built into the network architecture (cf. [7,8]). OBAI reconciles these two lines of research, by extending object-based visual inference models with action-based dynamics. We showed that OBAI can accurately track and predict the dynamics (and other properties) of simple objects whose movements are perturbed by actions – an important prerequisite for an agent to plan its own actions. In addition, we presented an efficient method for internalizing goals as a preference distribution over latent states, and showed that the agent can infer the actions necessary to accomplish these goals, at the

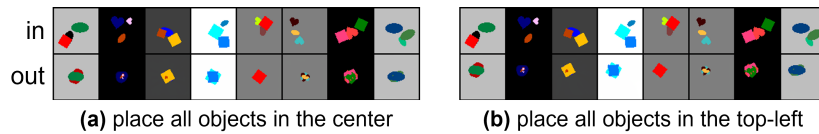


Figure 4. Goal-directed action planning. We give the network an arbitrary input image with three objects (in) and it infers the action that will move the state towards the learned preference, and imagines the resulting image (out). In (a), $\tilde{p}(s_{\text{true}})$ was biased towards the center of the image; in (b), it was biased more towards the top-left.

same abstract level of reasoning. While our results are preliminary, they are an important proof-of-concept, establishing the potential of our approach. In future work, we aim to scale OBAI to more naturalistic environments, and more cognitively demanding tasks.

5 Acknowledgements

RSvB is supported by Human Brain Project Specific Grant Agreement 3 grant ID 643945539: "SPIKEREFERENCE".

References

1. Lanillos, P., Dean-Leon, E., Cheng, G.: Yielding self-perception in robots through sensorimotor contingencies. *IEEE Transactions on Cognitive and Developmental Systems* **9**(2), 100–112 (2016)
2. Kourtzi, Z., Connor, C.E.: Neural representations for object perception: structure, category, and adaptive coding. *Annual review of neuroscience* **34**, 45–67 (2011)
3. Peters, B., Kriegeskorte, N.: Capturing the objects of vision with neural networks. *Nature Human Behaviour* **5**, 1127–1144 (9 2021). <https://doi.org/10.1038/s41562-021-01194-6>
4. Greff, K., Kaufman, R.L., Kabra, R., Watters, N., Burgess, C., Zoran, D., Matthey, L., Botvinick, M., Lerchner, A.: Multi-object representation learning with iterative variational inference. In: *International Conference on Machine Learning*. pp. 2424–2433. PMLR (2019)
5. Locatello, F., Weissenborn, D., Unterthiner, T., Mahendran, A., Heigold, G., Uszkoreit, J., Dosovitskiy, A., Kipf, T.: Object-centric learning with slot attention. *Advances in Neural Information Processing Systems* **33**, 11525–11538 (2020)
6. Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., et al.: Human-level control through deep reinforcement learning. *Nature* **518**(7540), 529–533 (2015)
7. Veerapaneni, R., Co-Reyes, J.D., Chang, M., Janner, M., Finn, C., Wu, J., Tenenbaum, J.B., Levine, S.: Entity abstraction in visual model-based reinforcement learning (10 2019), <http://arxiv.org/abs/1910.12827>
8. Watters, N., Matthey, L., Bosnjak, M., Burgess, C.P., Lerchner, A.: Cobra: Data-efficient model-based rl through unsupervised object discovery and curiosity-driven exploration (5 2019), <http://arxiv.org/abs/1905.09275>
9. Parr, T., Pezzulo, G., Friston, K.J.: *Active inference: the free energy principle in mind, brain, and behavior*. MIT Press (2022)
10. Lanillos, P., Meo, C., Pezzato, C., Meera, A.A., Baioumy, M., Ohata, W., Tschantz, A., Millidge, B., Wisse, M., Buckley, C.L., et al.: Active inference in robotics and artificial agents: Survey and challenges. *arXiv preprint arXiv:2112.01871* (2021)
11. Zablotskaia, P., Dominici, E.A., Sigal, L., Lehmman, A.M.: Unsupervised video decomposition using spatio-temporal iterative inference (2020), <http://arxiv.org/abs/2006.14727>
12. Marino, J., Yue, Y., Mandt, S.: Iterative amortized inference. *35th International Conference on Machine Learning, ICML 2018* **8**, 5444–5462 (7 2018), <http://arxiv.org/abs/1807.09356>

13. Matthey, L., Higgins, I., Hassabis, D., Lerchner, A.: dsprites: Disentanglement testing sprites dataset. <https://github.com/deepmind/dsprites-dataset/> (2017)
14. Active inference and epistemic value. *Cognitive Neuroscience* **6**, 187–214 (2015). <https://doi.org/10.1080/17588928.2015.1020053>
15. Sajid, N., Ball, P.J., Parr, T., Friston, K.J.: Active inference: Demystified and compared. *Neural Computation* **33**, 674–712 (3 2021). https://doi.org/10.1162/neco_a_01357, <https://direct.mit.edu/neco/article/33/3/674-712/97486>
16. Friston, K.: A free energy principle for a particular physics (6 2019), <http://arxiv.org/abs/1906.10184>
17. Da Costa, L., Parr, T., Sajid, N., Veselic, S., Neacsu, V., Friston, K.: Active inference on discrete state-spaces: a synthesis. *Journal of Mathematical Psychology* **99**, 102447 (2020)
18. Millidge, B., Tschantz, A., Buckley, C.L.: Whence the expected free energy? (4 2020), <http://arxiv.org/abs/2004.08128>
19. Jang, E., Gu, S., Poole, B.: Categorical reparameterization with gumbel-softmax (11 2016), <http://arxiv.org/abs/1611.01144>

Appendix

A Active-dSprites

Active-dSprites can be thought of as an "activated" version of the various multi-dSprites datasets that have been used in previous work on object-based visual inference (e.g. [4,5]). Not only does it include dynamics, but these dynamics can be acted on by an agent. Thus, active-dSprites is an interactive environment, rather than a dataset.

Objects in the active-dSprites environment are 2.5-D shapes (squares, ellipses and hearts): they have no depth dimension of their own, but can occlude each other within the depth dimension of the image. When an active-dSprites instance is initialized, object shapes, positions, sizes and colors are all sampled Uniformly at random. Initial velocities are drawn from a Normal distribution with mean 0 and standard deviation 4 (in units of pixels). Shape colors are sampled at discrete intervals spanning the full range of RGB-colors. Shapes are presented in random depth order against a solid background with a random grayscale color. Accelerations in action fields (at those locations that have been selected to incur a non-zero acceleration) are drawn from a Normal distribution with mean 0 and s.d. of 4.

B Network architectures

The OBAI architecture discussed in this paper consists of two separate IAI modules, each of which in turn contains a refinement and a decoder module. The first IAI module concerns the inference of the state beliefs $q(\{\mathbf{s}^{\dagger(k)}\})$ – we term this the *state inference module*. The second IAI module infers the object action beliefs $q(\{\mathbf{a}^{(k)}\})$, and we refer to this as the *action inference module*.

We ran inference for a total of $F \times 4$ iterations, where F is the number of frames in the input. Inference initially concerns just the first video frame, and beliefs for this frame are initialized to λ_0 , which is learned during training. After every 4 iterations, an additional frame is added to the inference window, and the beliefs for this new frame are initialized predictively, by extrapolating the object dynamics inferred up to that point. This procedure minimizes the risk that object representations are swapped between slots across frames, which can constitute a local minimum for the ELBO loss and leads to poor inference. We trained all IAI modules with $K=4$ object slots.

B.1 State inference module

This module used a latent dimension of 16. Note that, in the output of the refinement network, this number is doubled once as each latent belief is encoded by a mean and variance, and then doubled again as we represent (and infer) both the states and their first-order derivatives. In the decoder, the latent dimension is doubled only once, as the state derivatives do not enter into the reconstruction

of a video frame. As in IODINE [4], we use a spatial broadcast decoder, meaning that the latent beliefs are copied along a spatial grid with the same dimensions as a video frame, and each latent vector is concatenated with the (x, y) coordinate of its grid location, before passing through a stack of transposed convolution layers. Decoder and refinement network architectures are summarized in the tables below. The refinement network takes in 16 image-sized inputs, which are identical to those used in IODONE [4], except that we omit the leave-one-out likelihoods. Vector-sized inputs join the network after the convolutional stage (which processes only the image-sized inputs), and consist of the variational parameters and (stochastic estimates of) their gradients.

Decoder

Type	Size/#Chan.	Act. func.	Comment
Input (λ)	32		
Broadcast	34		Appends coordinate channels
Conv ^T 5 × 5	32	ELU	
Conv ^T 5 × 5	32	ELU	
Conv ^T 5 × 5	32	ELU	
Conv ^T 5 × 5	32	ELU	
Conv ^T 5 × 5	4		Outputs RGB + mask

Refinement network

Type	Size/#Chan.	Act. func.	Comment
Linear	64		
LSTM	128	tanh	
Concat [..., λ , $\nabla_{\lambda}\mathcal{L}$]	256		Appends vector-sized inputs
Linear	128	ELU	
Flatten	800		
Conv 5 × 5	32	ELU	
Conv 5 × 5	32	ELU	
Conv 5 × 5	32	ELU	
Inputs	16		

B.2 Action inference module

The action inference module does not incorporate a decoder network, as the quality of the action beliefs is computed by evaluating equation 7 and plugging this into the ELBO loss from equation 8. While this requires some additional tricks (see Appendix D), no neural network is required for this. This module does include a (shallow) refinement network, which is summarized in the table below.

This network takes as input the current variational parameters $\lambda_{\mathbf{a}}^{(k)}$ (2 means and 2 variances), their gradients, and the ‘expected object action’, $\sum_i \hat{m}_{ik} \psi_i$.

Refinement network

Type	Size/#Chan.	Act. func.	Comment
Linear	4		
LSTM	32	tanh	
Inputs	10		

C Training procedure

The above network architecture was trained on pre-generated experience with the active-dSprites environment, as described in the main text. The training set comprised 50,000 videos of 4 frames each. An additional validation set of 10,000 videos was constructed using the same environment parameters as the training set, but using a different random seed. Training was performed using the ADAM optimizer [REF] with default parameters and an initial learning rate of 3×10^{-4} . This learning rate was reduced automatically by a factor 3 whenever the validation loss had not decreased in the last 10 training epochs, down to a minimum learning rate of 3×10^{-5} . Training was performed with a batch size of 64 (16×4 GPUs), and was deemed to have converged after 245 epochs.

C.1 Modified ELBO loss

OBAI optimizes an ELBO loss for both learning and inference. The basic form of this loss is given by equation 8. In practice, we modify this loss in two ways (similar to previous work, e.g. [4]). First, we re-weight the reconstruction term in the ELBO loss as follows:

$$\mathcal{L}_\beta = - \sum_{t=0}^T \left[\mathcal{H} \left(q \left(\{\mathbf{s}_t^{\dagger(k)}, \mathbf{a}_t^{(k)}\} \right) \right) + \beta E_{q(\{\mathbf{s}_t^{(k)}\})} [\log p(\mathbf{o}_t | \{\mathbf{s}_t^{(k)}\})] \right. \\ \left. + \sum_k E_{q(\mathbf{a}_t^{(k)})} [\log p(\mathbf{a}_t^{(k)} | \Psi_t)] + \sum_k E_{q(\mathbf{s}_t^{\dagger(k)}, \mathbf{s}_{t-1}^{\dagger(k)}, \mathbf{a}_{t-1}^{(k)})} [\log p(\mathbf{s}_t^{\dagger(k)} | \mathbf{s}_{t-1}^{\dagger(k)}, \mathbf{a}_{t-1}^{(k)})] \right] \quad (15)$$

Second, we train the network to minimize not just the loss at the end of the inference iterations through the network, but a composite loss that also includes the loss after earlier iterations. Let $\mathcal{L}_\beta^{(n)}$ be the loss after n inference iterations, then the composite loss is given by:

$$\mathcal{L}_{\text{comp}} = \sum_{n=1}^{N_{\text{iter}}} \frac{n}{N_{\text{iter}}} \mathcal{L}_\beta^{(n)} \quad (16)$$

C.2 Hyperparameters

OBAI includes a total of 4 hyperparameters: (1) the loss-reweighting coefficient β (see above); (2) the variance of the pixels around their predicted values, σ_o^2 ; (3) the variance of the noise in the latent space dynamics, σ_s^2 ; and (4) the variance of the noise in the object actions, σ_ψ^2 . The results described in the current work were achieved with the following settings:

Param.	Value
β	5.0
σ_o	0.3
σ_s	0.1
σ_ψ	0.3

D Computing $E_{q(\mathbf{a}^{(k)})}[\log p(\mathbf{a}^{(k)}|\Psi)]$

The expectation under $q(\mathbf{a}^{(k)})$ of $\log p(\mathbf{a}^{(k)}|\Psi)$, which appears in the ELBO loss (eq. 8), cannot be computed in closed form, because the latter log probability requires us to marginalize over all possible configurations of the pixel-to-object assignments, and to do so inside of the logarithm. That is:

$$\log p(\mathbf{a}^{(k)}|\Psi) = \sum_{\mathbf{m}} \log \left(p(\mathbf{a}^{(k)}|\Psi, \mathbf{m}) p(\mathbf{m}|\{\mathbf{s}^{(k)}\}) \right) \quad (17)$$

$$= \log \left(E_{p(\mathbf{m}|\{\mathbf{s}^{(k)}\})} [p(\mathbf{a}^{(k)}|\Psi, \mathbf{m})] \right) \quad (18)$$

However, note that within the ELBO loss, we want to maximize the expected value of this quantity (as its negative appears in the ELBO, which we want to minimize). From Jensen’s inequality, we have:

$$E_{p(\mathbf{m}|\{\mathbf{s}^{(k)}\})} [\log p(\mathbf{a}^{(k)}|\Psi, \mathbf{m})] \leq \log \left(E_{p(\mathbf{m}|\{\mathbf{s}^{(k)}\})} [p(\mathbf{a}^{(k)}|\Psi, \mathbf{m})] \right) \quad (19)$$

Therefore, the l.h.s. of this equation provides a lower bound on the quantity we want to maximize. Thus, we can approximate our goal by maximizing this lower bound instead. This is convenient, because this lower bound, and its expectation under $q(\mathbf{a}^{(k)})$ can be approximated through sampling:

$$E_{q(\mathbf{a}^{(k)})} \left[E_{p(\mathbf{m}|\{\mathbf{s}^{(k)}\})} [\log p(\mathbf{a}^{(k)}|\Psi, \mathbf{m})] \right] \approx \frac{1}{N_{\text{samples}}} \sum_j \log p(\mathbf{a}_j^{(k)*}|\Psi, \mathbf{m}_j^*) \quad (20)$$

$$= \frac{1}{N_{\text{samples}}} \sum_j \log \mathcal{N} \left(\mathbf{a}_j^{(k)*}; \sum_i \hat{m}_{jk}^{*(i)} \boldsymbol{\psi}_i, \sigma_\psi^2 \mathbf{I} \right) \quad (21)$$

$$\hat{m}_j^{*(i)} \sim p(m_i|\{\mathbf{s}^{(k)}\}), \quad \mathbf{a}_j^{(k)*} \sim q(\mathbf{a}^{(k)}), \quad \mathbf{s}_j^{(k)*} \sim q(\mathbf{s}^{(k)}) \quad (22)$$

where we slightly abuse notation in the sampling of the pixel assignments, as a vector is sampled from a distribution over a categorical variable. The reason

this results in a vector is because this sampling step uses the Gumbel-Softmax trick [19], which is a differentiable method for sampling categorical variables as "approximately one-hot" vectors. Thus, for every pixel i , we sample a vector $\hat{\mathbf{m}}_j^{*(i)}$, such that the k -th entry of this vector, $\hat{m}_{jk}^{*(i)}$, denotes the "soft-binary" condition of whether pixel i belongs to object k . In practice, we use $N_{\text{samples}} = 1$, based on the intuition that this will still yield a good approximation over many training instances, and that we rely on the refinement network to learn to infer good beliefs. The Gumbel-Softmax sampling method depends on a temperature τ , which we gradually reduce across training epochs, so that the samples gradually better approximate the ideal one-hot vectors.

It is worth noting that, as the entropy of $p(\mathbf{m}|\{\mathbf{s}^{(k)}\})$ decreases (i.e. as object slots "become more certain" about which pixels are theirs), the bound in equation 19 becomes tighter. In the limit as the entropy becomes 0, the network is perfectly certain about the pixel assignments, and so the distribution collapses to a point mass. The expectation then becomes trivial, and so the two sides of eq. 19 become equal. Sampling the pixel assignments is equally trivial in this case, as the distribution has collapsed to permit only a single value for each assignment. In short, at this extreme point, the procedure becomes entirely deterministic. In our data, we typically observe very low entropy for $p(\mathbf{m}|\{\mathbf{s}^{(k)}\})$, and so we likely operate in a regime close to the deterministic one, where the approximation is very accurate.