
Unsupervised Learning of Object Keypoints for Perception and Control

Tejas Kulkarni^{*1}, Ankush Gupta^{*1}, Catalin Ionescu¹, Sebastian Borgeaud¹,
Malcolm Reynolds¹, Andrew Zisserman^{1,2}, and Volodymyr Mnih¹

* indicates equal contribution

¹DeepMind, London

²VGG, Department of Engineering Science, University of Oxford

{tkulkarni, ankushgupta, cdi, sborgeaud, mareynolds, zisserman, vmnih}@google.com

Abstract

The study of object representations in computer vision has primarily focused on developing representations that are useful for image classification, object detection, or semantic segmentation as downstream tasks. In this work we aim to learn object representations that are useful for control and reinforcement learning (RL). To this end, we introduce *Transporter*, a neural network architecture for discovering concise geometric object representations in terms of *keypoints* or image-space coordinates. Our method learns from raw video frames in a fully unsupervised manner, by *transporting* learnt image features between video frames using a keypoint bottleneck. The discovered keypoints track objects and object parts across long time-horizons more accurately than recent similar methods. Furthermore, consistent long-term tracking enables two notable results in control domains – (1) using the keypoint co-ordinates and corresponding image features as inputs enables highly sample-efficient reinforcement learning; (2) learning to explore by controlling keypoint locations drastically reduces the search space, enabling deep exploration (leading to states unreachable through random action exploration) without any extrinsic rewards. Code for the model is available at: <https://github.com/deepmind/deepmind-research/tree/master/transporter>.

1 Introduction

End-to-end learning of feature representations has led to advances in image classification [21], generative modeling of images [8] and agents which outperform expert humans at game play [26, 33]. However, this training procedure induces task-specific representations, especially in the case of reinforcement learning, making it difficult to re-purpose the learned knowledge for future unseen tasks. On the other hand, humans explicitly learn notions of objects, relations, geometry and cardinality in a task-agnostic manner [34] and re-purpose this knowledge to future tasks. There has been extensive research inspired by psychology and cognitive science on explicitly learning object-centric representations from pixels. Both instance and semantic segmentation has been approached using supervised [25, 28] and unsupervised learning [3, 10, 16, 11, 18, 24, 7] methods. However, the representations learned by these methods do not explicitly encode fine-grained locations and orientations of object parts, and thus they have not been extensively used in the control and reinforcement learning literature. We argue that being able to precisely control objects and object parts is at the root of many complex sensory motor behaviours.

In recent work, object keypoint or landmark discovery methods [42, 17] have been proposed to learn representations that precisely represent locations of object parts. These methods predict a set of Cartesian co-ordinates of keypoints denoting the salient locations of objects given image frame(s).

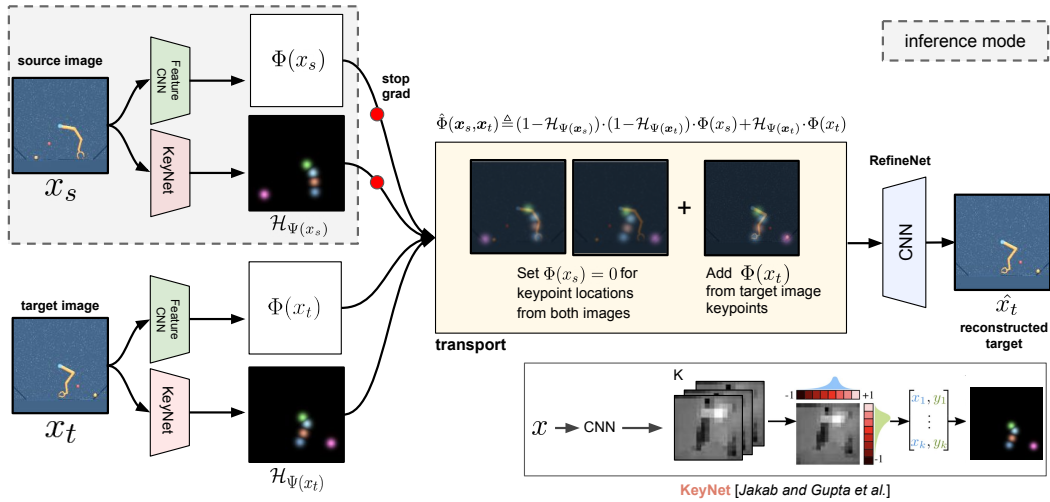


Figure 1: **Transporter**. Our model leverages object motion to discover keypoints by learning to transform a source video frame (x_s) into another target frame (x_t) by *transporting* image features at the discovered object locations. During training, spatial feature maps $\Phi(x)$ and keypoints co-ordinates $\Psi(x)$ are predicted for both the frames using a ConvNet and the fully-differentiable *KeyNet* [17] respectively. The keypoint co-ordinates are transformed into Gaussian heatmaps (same spatial dimensions as feature maps) $\mathcal{H}_{\Psi(x)}$. We perform two operations in the transport phase: (1) the features of the source frame are set to zero at both locations $\mathcal{H}_{\Psi(x_s)}$ and $\mathcal{H}_{\Psi(x_t)}$; (2) the features in the source image $\Phi(x_s)$ at the target positions $\Psi(x_t)$ are replaced with the features from the target image $\mathcal{H}_{\Psi(x_t)} \cdot \Phi(x_t)$. The final refinement ConvNet (which maps from the transported feature map to an image) then has two tasks: (i) to inpaint the missing features at the source position; and (ii) to clean up the image around the target positions. During inference, keypoints can be extracted for a *single* frame via a feed-forward pass through the *KeyNet* (Ψ).

However, as we will show, the existing methods struggle to accurately track keypoints under the variability in number, size, and motion of objects present in common RL domains.

We propose *Transporter*, a novel architecture to explicitly discover spatially, temporally and geometrically aligned keypoints given only videos. After training, each keypoint represents and tracks the co-ordinate of an object or object part even as it undergoes deformations (see fig. 1 for illustrations). As we will show, *Transporter* learns more accurate and more consistent keypoints on standard RL domains than existing methods. We will then showcase two ways in which the learned keypoints can be used for control and reinforcement learning. First, we show that using keypoints as inputs to policies instead of RGB observations leads to drastically more data efficient reinforcement learning on Atari games. Second, we show that by learning to control the Cartesian coordinates of the keypoints in the image plane we are able to learn skills or options [35] grounded in pixel observations, which is an important problem in reinforcement learning. We evaluate the learned skills by using them for exploration and show that they lead to much better exploration than primitive actions, especially on sparse reward tasks. Crucially, the learned skills are task-agnostic because they are learned without access to any rewards.

In summary, our key contributions are:

- *Transporter* learns state of the art object keypoints across a variety of commonly used RL environments. Our proposed architecture is robust to varying number, size and motion of objects.
- Using learned keypoints as state input leads to policies that perform better than state-of-the-art model-free and model-based reinforcement learning methods on several Atari environments, while using only up to 100k environment interactions.
- Learning skills to manipulate the most controllable keypoints provides an efficient action space for exploration. We demonstrate drastic reductions in the search complexity for exploring challenging Atari environments. Surprisingly, our action space enables random agents to play several Atari games without rewards and any task-dependent learning.

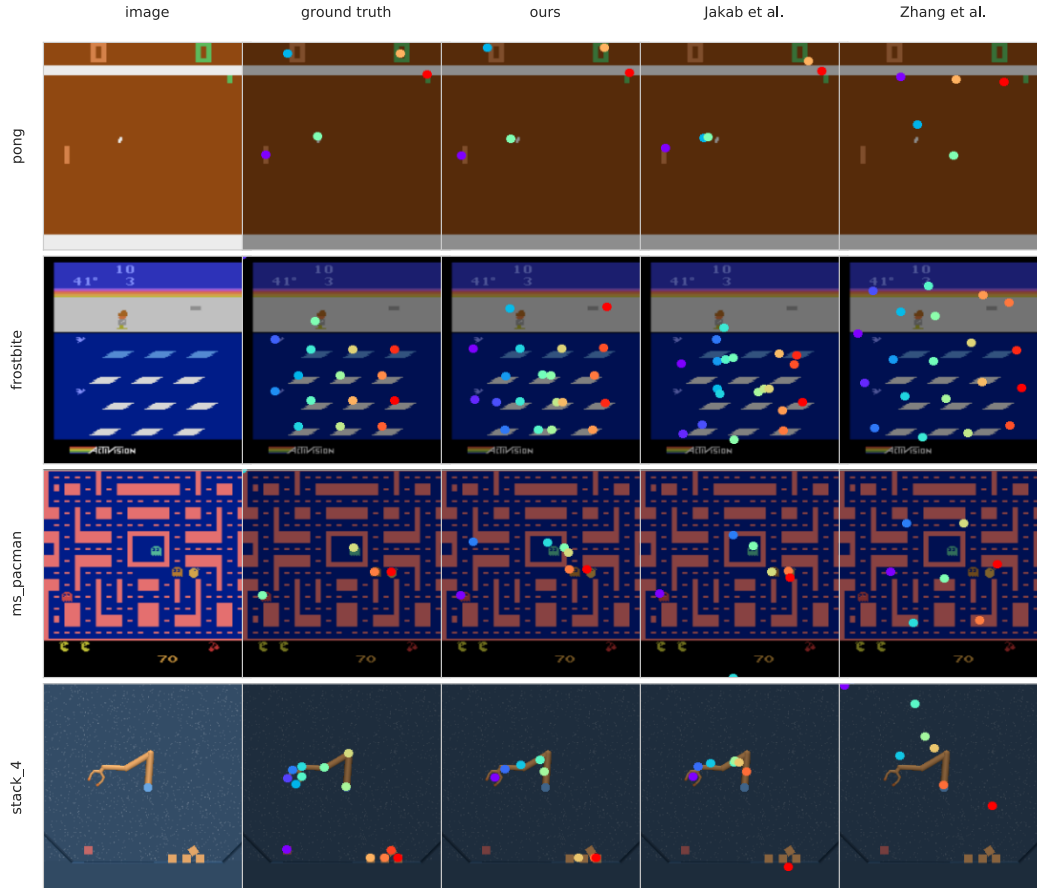


Figure 2: **Keypoint visualisation.** Visualisations from our and state-of-the-art unsupervised object keypoint discovery methods: Jakab *et al.* [17] and Zhang *et al.* [42] on Atari ALE [1] and Manipulator [37] domains. Our method learns more spatially aligned keypoints, *e.g.* `frostbite` and `stack_4` (see section 4.1). Quantitative evaluations are given in fig. 4 and further visualisations in the supplementary material.

2 Related Work

Our work is related to the recently proposed literature on unsupervised object keypoint discovery [42, 17]. Most notably, Jakab and Gupta *et al.* [17] proposed an encoder-decoder architecture with a differentiable keypoint bottleneck. We reuse their bottleneck architecture but add a crucial new inductive bias – the feature transport mechanism – to constrain the representation to be more spatially aligned compared to all baselines. The approach in Zhang *et al.* [42] discovers keypoints using single images and requires privileged information about temporal transformations between frames in form of optical flow. This approach also requires multiple loss and regularization terms to converge. In contrast, our approach does not require access to these transformations and learns keypoints with a simple pixel-wise L2 loss function. Other works similarly either require known transformations or output dense correspondences instead of discrete landmarks [38, 32, 36, 40]. Deep generative models with structured bottlenecks have recently seen a lot of advances [4, 22, 39, 41, 13] but they do not explicitly reason about geometry.

Unsupervised learning of object keypoints has not been widely explored in the control literature, with the notable exception of [6]. However, this model uses a full-connected layer for reconstruction and therefore can learn non-spatial latent embeddings similar to a baseline we consider [17]. Moreover, similar to [42] their auto-encoder reconstructs *single* frames and hence does not learn to factorize geometry. Object-centric representations have also been studied in the context of intrinsic motivation, hierarchical reinforcement learning and exploration. However, existing approaches either require hand-crafted object representations [23] or have not been shown to capture fine-grained representations over long temporal horizons [16].

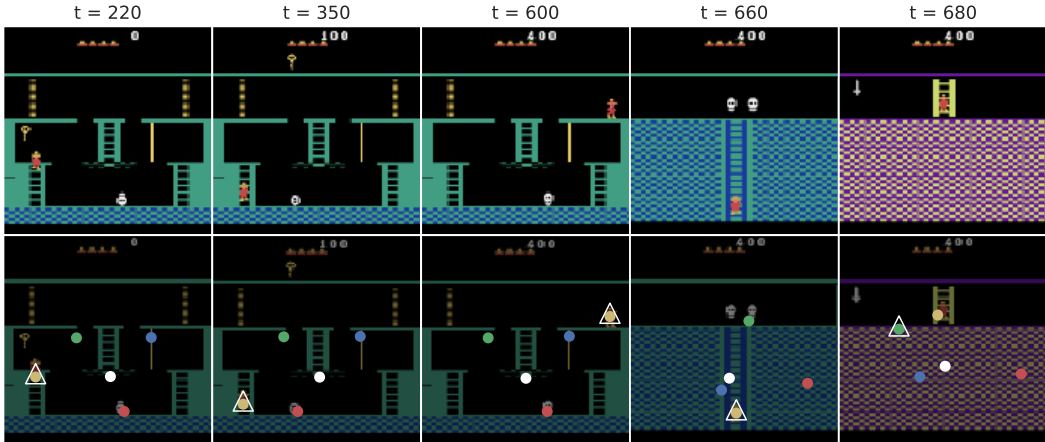


Figure 3: **Temporal consistency of keypoints.** Our learned keypoints are temporally consistent across hundreds of environment steps, as demonstrated in this classical hard exploration game called montezuma’s revenge [1]. Additionally, we also predict the most controllable keypoint denoted by the triangular markers, without using any environment rewards. This prediction often corresponds to the avatar in the game and it is consistently tracked across different parts of the state space. See section 4.2.2 for further discussion.

3 Method

In section 3.1 we first detail our model for unsupervised discovery of object keypoints from videos. Next, we describe the application of the learned object keypoints to control for – (1) data-efficient reinforcement learning (section 3.2.1) and (2) learning keypoint based options for efficient exploration (section 3.2.2).

3.1 Feature Transport for learning Object Keypoints

Given an image \mathbf{x} , our objective is to extract K 2-dimensional image locations or *keypoints*, $\Psi(\mathbf{x}) \in \mathbb{R}^{K \times 2}$, which correspond to locations of objects or object-parts without any manual labels for locations. We follow the formulation of [17] and assume access to frame pairs $\mathbf{x}_s, \mathbf{x}_t$ collected from some trajectories such that the frames differ only in objects’ pose / geometry or appearance. The learning objective is to reconstruct the second frame \mathbf{x}_t from the first \mathbf{x}_s . This is achieved by computing ConvNet (CNN) feature maps $\Phi(\mathbf{x}_s), \Phi(\mathbf{x}_t) \in \mathbb{R}^{H' \times W' \times D}$ and extracting K 2D locations $\Psi(\mathbf{x}_s), \Psi(\mathbf{x}_t) \in \mathbb{R}^{K \times 2}$ by marginalising the keypoint-detector feature-maps along the image dimensions (as proposed in [17]). A *transported* feature map $\hat{\Phi}(\mathbf{x}_s, \mathbf{x}_t)$ is generated by suppressing both sets of keypoint locations in $\Phi(\mathbf{x}_s)$ and compositing in the featuremaps around the keypoints from \mathbf{x}_t :

$$\hat{\Phi}(\mathbf{x}_s, \mathbf{x}_t) \triangleq (1 - \mathcal{H}_{\Psi(\mathbf{x}_s)}) \cdot (1 - \mathcal{H}_{\Psi(\mathbf{x}_t)}) \cdot \Phi(\mathbf{x}_s) + \mathcal{H}_{\Psi(\mathbf{x}_t)} \cdot \Phi(\mathbf{x}_t) \quad (1)$$

where \mathcal{H}_{Ψ} is a heatmap image containing fixed-variance isotropic Gaussians around each of the K points specified by Ψ . A final CNN with small-receptive field refines the transported reconstruction $\hat{\Phi}(\mathbf{x}_s, \mathbf{x}_t)$ to regress the target frame $\hat{\mathbf{x}}_t$. We use pixel-wise squared- ℓ_2 reconstruction error $\|\mathbf{x}_t - \hat{\mathbf{x}}_t\|_2^2$ for end-to-end learning. The keypoint network Ψ learns to track moving entities between frames to enable successful reconstruction.

In words, (i) the features in the source image $\Phi(\mathbf{x}_s)$ at the target positions $\Psi(\mathbf{x}_t)$ are replaced with the features from the target image $\mathcal{H}_{\Psi(\mathbf{x}_t)} \cdot \Phi(\mathbf{x}_t)$ — this is the *transportation*; and (ii) the features at the source position $\Psi(\mathbf{x}_s)$ are set to zero. The refine net (which maps from the transported feature map to an image) then has two tasks: (i) to inpaint the missing features at the source position; and (ii) to clean up the image around the target positions. Refer to fig. 1 for a concise description of our method.

Note, unlike [17] who regress the target frame from *stacked* target keypoint heatmaps $\mathcal{H}_{\Psi(\mathbf{x}_t)}$ and source image features $\Phi(\mathbf{x}_s)$, we enforce explicit *spatial transport* for stronger correlation with image locations leading to more robust long-term tracking (section 4.1).

3.2 Object Keypoints for Control

Consider a Markov Decision Process (MDP) $(\mathcal{X}, \mathcal{A}, \mathcal{T}, r)$ with pixel observations $\mathbf{x} \in \mathcal{X}$ as states, actions $a \in \mathcal{A}$, transition function $T: \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{X}$ and reward function $r: \mathcal{X} \rightarrow \mathbb{R}$. *Transporter* keypoints provide a concise visual state-abstraction which enables faster learning (section 3.2.1). Further, in 3.2.2 we demonstrate task-independent exploration by learning to control the keypoint coordinates.

3.2.1 Data-efficient reinforcement learning

Our first hypothesis is that task-agnostic learning of object keypoints can enable fast learning of a policy. This is because once keypoints are learnt, the control policy can be much simpler and does not have to relearn visual features using temporal difference learning. In order to test this hypothesis, we use a variant of the neural fitted Q-learning framework [30] with learned keypoints as input and a recurrent neural network Q function to output behaviors. Specifically, the agent observes the *thermometer encoded* [2] keypoint coordinates $\Psi(\mathbf{x}_t)$, and also the features $\Phi(\mathbf{x}_t)$ under the keypoint locations obtained by spatially averaging the feature tensor (Φ) multiplied with (Gaussian) heat-maps (\mathcal{H}_Ψ). *Transporter* is pre-trained by collecting data using a random policy and without any reward functions (see supplementary material for details). Then, *Transporter* network weights are fixed during behavior learning from environment rewards.

3.2.2 Keypoint-based options for efficient exploration

Our second hypothesis is that learned keypoints can enable significantly better task-independent exploration. Typically, raw actions are randomly sampled to bootstrap goal-directed policy learning. This exploration strategy is notoriously inefficient. We leverage the *Transporter* representation to learn a new action space. The actions are now skills grounded in the control of co-ordinate values of each keypoint. This idea has been explored in the reinforcement learning community [23, 16] but it has been hard to learn spatial features with long temporal consistency. Here we show that *Transporter* is particularly amenable to this task. We show that randomly exploring in this space leads to significantly more rewards compared to raw actions. Our learned action space is agnostic to the control algorithm and hence other exploration algorithms [27, 5, 29] can also benefit from using it.

To do this, we define $K \times 4$ intrinsic reward functions using the keypoint locations, similar to the *VisEnt* agent [16]. Each reward function corresponds to how much each keypoint moves in the 4 cardinal directions (up, down, left, right) between consecutive observations. We learn a set of $K \times 4$ Q functions $\{Q_{i,j} | i \in \{1, \dots, K\}, j \in \{1, 2, 3, 4\}\}$ to maximise each of the following reward functions: $r_{i,1} = \Psi_x^i(\mathbf{x}_{t+1}) - \Psi_x^i(\mathbf{x}_t)$, $r_{i,2} = \Psi_x^i(\mathbf{x}_t) - \Psi_x^i(\mathbf{x}_{t+1})$, $r_{i,3} = \Psi_y^i(\mathbf{x}_{t+1}) - \Psi_y^i(\mathbf{x}_t)$, $r_{i,4} = \Psi_y^i(\mathbf{x}_t) - \Psi_y^i(\mathbf{x}_{t+1})$. These functions correspond to increasing/decreasing the x and y coordinates respectively. The Q functions are trained using n-step $Q(\lambda)$.

During training, we randomly sample a particular Q function to act with and commit to this choice for T timesteps before resampling. All Q functions are trained using experiences generated from all policies via a shared replay buffer. Randomly exploring in this Q space can already reduce the search space as compared to raw actions. During evaluation, we further reduce this search space using a fixed *controllability policy* $\pi_{Q_{\text{gap}}}$ to select the single ‘‘most controllable’’ keypoint, where

$$\pi_{Q_{\text{gap}}}(s) = \underset{i}{\operatorname{argmax}} \frac{1}{4} \sum_{j=1}^4 \max_a Q_{i,j}(s;a) - \min_a Q_{i,j}(s;a). \quad (2)$$

$\pi_{Q_{\text{gap}}}$ picks keypoints for which actions lead to more prospective change in all spatial directions than all other keypoints. For instance, in Atari games this corresponds to the avatar which is directly controllable on the screen. Our random exploration policy commits to the $Q_{i,j}$ function corresponding to the keypoint i selected as above and a direction j sampled uniformly at random for T timesteps and then resamples. Consider a sequence of 100 actions with 18 choices before receiving rewards, which is typical in hard exploration Atari games (e.g. *montezuma’s revenge*). A random action agent would need to search in the space of 18^{100} raw actions. However, observing 5 keypoints and $T = 20$ only has $(5 \times 4)^{100/20}$, giving a search space reduction of 10^{100} . The search space reduces further when we explore with the most controllable keypoints. Since our learned action space is agnostic to the control mechanism, we evaluate them by randomly searching in this space versus raw actions. We measure extrinsically defined game score as the metric to evaluate the effectiveness of both search procedures.

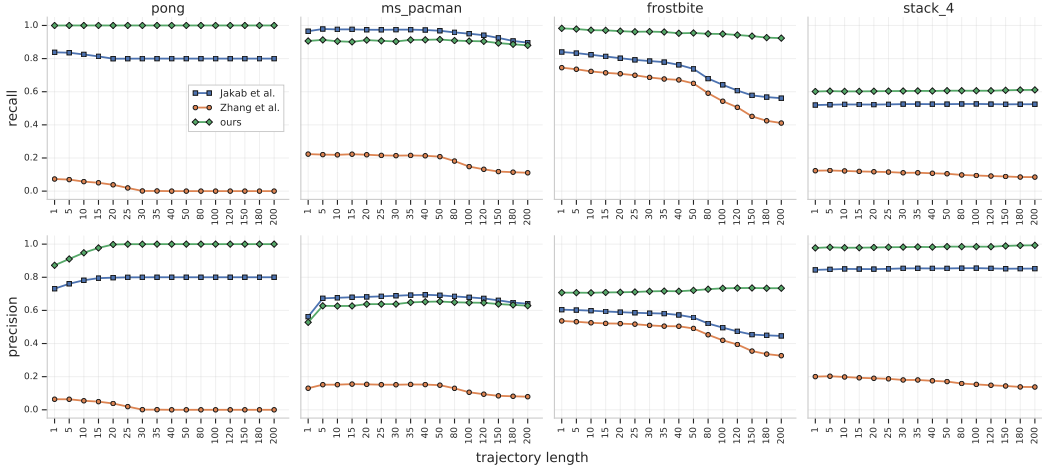


Figure 4: **Long-term tracking evaluation.** We compare long-term tracking ability of our keypoint detector against Jakab *et al.* [17] and Zhang *et al.* [42] (visualisations in fig. 2 and supplementary material). We report precision and recall for trajectories of varying lengths (lengths = 1 – 200 frames; each frame corresponds to 4 action repeats) against ground-truth keypoints on Atari ALE [1] and Manipulator [37] domains. Our method significantly outperforms the baselines on all games (100% on pong), except for ms_pacman where we perform similarly especially for long trajectories (length = 200). See section 4.1 for further discussion.

4 Experiments

In section 4.1 we first evaluate the long-term tracking ability of our object keypoint detector. Next, in section 4.2 we evaluate the application of the keypoint detector on two control tasks — comparison against state-of-the-art model-based and model-free methods for data-efficient learning on Atari ALE games [1] in section 4.2.1, and in section 4.2.2 examine efficient exploration by learning to control the discovered keypoints; we demonstrate reaching states otherwise unreachable through random explorations on raw-actions, and also recover the agent *self* as the most-controllable keypoint. For implementation details, please refer to the supplementary material.

Datasets. We evaluate our method on Atari ALE [1] and Manipulator [37] domains. We chose representative levels with large variations in the type and number of objects. (1) For evaluating long-term tracking of object keypoints section 4.1 we use — pong, frostbite, ms_pacman, and stack_4 (manipulator with blocks). (2) For data-efficient reinforcement learning (section 4.2.1) we train on diverse data collected using random exploration on the Atari games indicated in fig. 6. (3) For keypoints based efficient-exploration (section 4.2.2) we evaluate on one of the most difficult exploration game — montezuma_revenge, along with ms_pacman and seaquest.

A random policy executes actions and we collect a trajectory of images before the environment resets; details for data generation are presented in the supplementary material. We sample the source and target frames x_s, x_t randomly within a temporal offset of 1 to 20 frames, corresponding to small or significant changes in the the configuration between these two frames respectively. For Atari ground-truth object locations are extracted from the emulated RAM using hand crafted per-game rules and for Manipulator it is extracted from the simulator geoms. The number of keypoints K is set to the maximum number of moving entities in each environment.

4.1 Evaluating Object Keypoint Predictions

Baselines. We compare our method against state-of-the-art methods for unsupervised discovery of object landmarks — (1) Jakab *et al.* [17] and (2) Zhang *et al.* [42]. For (1) we use exactly the same architecture for Φ and Ψ as ours; for (2) we use the implementation released online by the authors where the image-size is set to 80×80 pixels. We train all the methods for 10^6 optimization steps and pick the best model checkpoint based on a validation set.

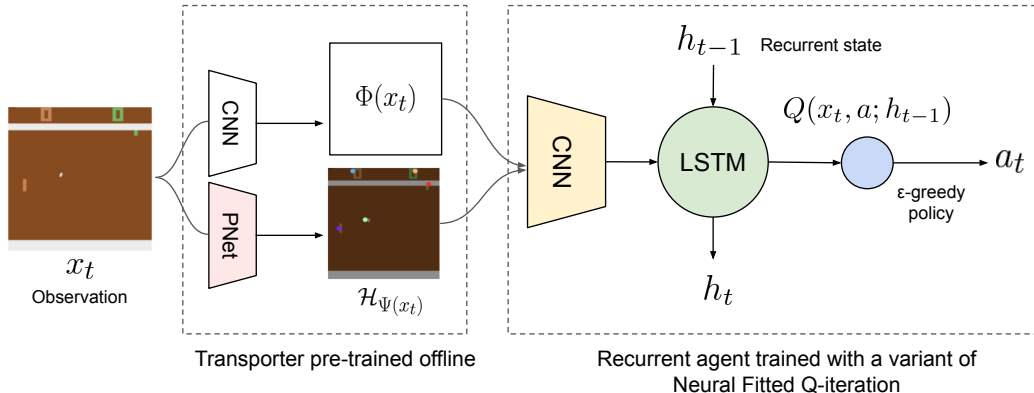


Figure 5: **Agent architecture for data-efficient reinforcement learning.** *Transporter* is trained off-line with data collected using a random policy. A recurrent variant of the neural-fitted Q-learning algorithm [30] rapidly learns control policies using keypoint co-ordinates and features at the corresponding locations given game rewards.

Game	KeyQN (ours)	SimPLe	Rainbow	PPO (100k)	Human	Random
breakout	19.3 (4.5)	12.7 (3.8)	3.3 (0.1)	5.9 (3.3)	31.8	1.7
frostbite	388.3 (142.1)	254.7 (4.9)	140.1 (2.7)	174.0 (40.7)	4334.7	65.2
ms_pacman	999.4 (145.4)	762.8 (331.5)	364.3 (20.4)	496.0 (379.8)	15693.0	307.3
pong	10.8 (5.7)	5.2 (9.7)	-19.5 (0.2)	-20.5 (0.6)	9.3	-20.7
seaquest	236.7 (22.2)	370.9 (128.2)	206.3 (17.1)	370.0 (103.3)	20182.0	-20.7

Figure 6: **Atari Mean Scores.** Mean scores (and std-dev in parentheses) obtained by our method (three random seeds) in comparison with Rainbow [12], SimPLe [19] and PPO [31] trained on 100K steps (400K frames). See section 4.2.1 for details. Numbers (except for KeyQN) taken from [19].

Metrics. We measure the precision and recall of the detected keypoint trajectories, varying their lengths from 1 to 200 frames (200 frames \approx 13 seconds @ 15-fps with action-repeat of 4) to evaluate long-term consistency of the keypoint detections crucial for control. The average Euclidean distance between each detected and ground-truth trajectory is computed. The time-steps where a ground-truth object is absent are ignored in the distance computation. Distances above a threshold (ϵ) are excluded as potential matches.¹ One-to-one assignments between the trajectories are then computed using min-cost linear sum assignment, and the matches are used for reporting precision and recall.

Results. Figure 2 visualises the detections while fig. 4 presents precision and recall for varying trajectory lengths. *Transporter* consistently tracks the salient object keypoints over long time horizons and outperforms the baseline methods on all environments, with the notable exception of [17] on pacman where our method is slightly worse but achieves similar performance for long-trajectories.

4.2 Using Keypoints for Control

4.2.1 Data-efficient Reinforcement Learning on Atari

We demonstrate that using the learned keypoints and corresponding features within a reinforcement learning context can lead to data-efficient learning in Atari games. Following [19], we trained our Keypoint Q-Network (KeyQN) architecture for 100,000 interactions, which corresponds to 400,000 frames. As shown in fig. 6, our approach is better than the state-of-the-art model-based SimPLe architecture [19], as well as the model-free Rainbow architecture [12] on four out of five games. Applying

¹The threshold value (ϵ) for evaluation is set to the average ground-truth spatial extent of entities for each environment (see supplementary material for details).

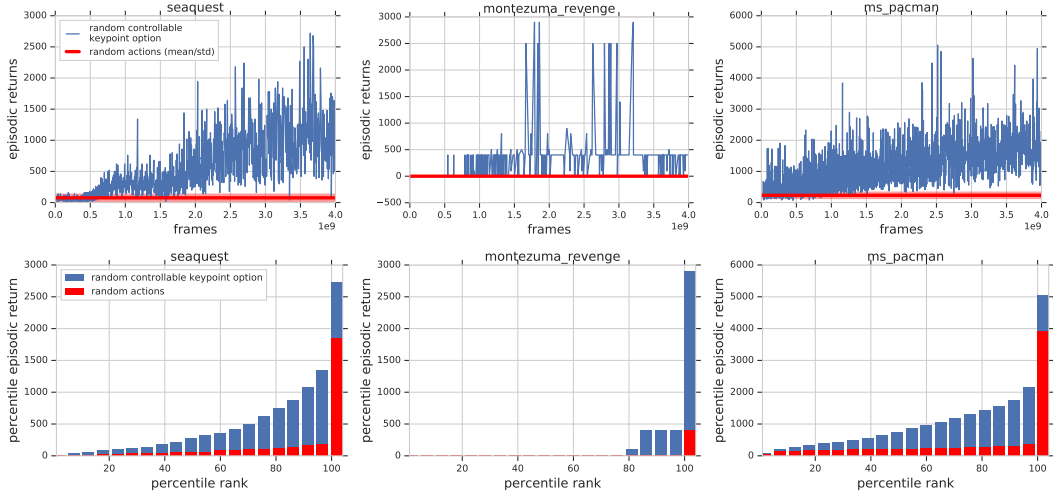


Figure 7: **Exploration using random actions versus random (most controllable) keypoint option / skills:** (*first row*) We perform random actions in the environment for all methods (without reward) and record the mean and standard deviation of episodic returns across 4 billion frames. With the same frame budget, we simultaneously learn the most controllable keypoint and randomly explore in the space of its co-ordinates (to move it *left, right, top, down*). The options model becomes better with training (using only intrinsic rewards) and this leads to higher extrinsically defined episodic returns. Surprisingly, our learned options model is able to play several Atari games via random sampling of options. This is possible by learning skills to move the discovered game avatar as far as possible without dying. (*second row*) We measure the percentile episodic return reached for all methods. Our approach outperforms the baseline, both in terms of efficient and robust exploration of rare and rewarding states.

this approach to all Atari games will require training *Transporter* inside the reinforcement learning loop because pre-training keypoints on data from a random policy is insufficient for games where new objects or screens can appear. However, these experiments provide evidence that the right visual abstractions and simple control algorithms can produce highly data efficient reinforcement learning algorithms.

4.2.2 Efficient Exploration with Keypoints

We learn options/skills for efficient exploration from the object keypoints. We use a distributed off-policy learner similar to [14] using 128 actors and 4 GPUs. The agent network has a standard CNN architecture [26] with an LSTM with 256 hidden units which feeds into a linear layer with $K \times 4 \times a$ units, where a is the number of actions. Our *Transporter* model is learnt simultaneously with all the control policies (no pre-training). We commit to the choice of Q function (corresponding to a keypoint and one of the four directions) for $T = 20$ steps (see section 3.2.2 for details). Actions are sampled using an ϵ -greedy random policy during training (ϵ is sampled from a log-uniform distribution over $[1e-4, 0.4]$), and greedily for evaluation. Quantitative results are shown in fig. 7. We also show qualitative results of the most controllable keypoint in fig. 3 and the supplementary material.

Our experiments clearly validate our hypothesis that using keypoints enables temporally extended exploration. As shown in fig. 7, our learned keypoint options consistently outperform the random actions baseline by a large margin. Encouragingly, our random options policy is able to play some Atari games by moving around the avatar (most controllable keypoint) in different parts of the state space without dying. For instance, the agent explores multiple rooms in Montezuma’s Revenge, a classical hard exploration environment in the reinforcement learning community. Similarly, our keypoint exploration learns to consistently move around the submarine in Seaquest and the avatar in Ms. Pacman. Most notably, this is achieved without rewards or (extrinsic) task-directed learning. Therefore our learned keypoints are stable enough to learn complex object-oriented skills in the Atari domain.

5 Conclusion

We demonstrate that it is possible to learn stable object keypoints across thousands of environment steps, without having access to task-specific reward functions. Therefore, object keypoints could provide a flexible and re-purposable representation for efficient control and reinforcement learning. Scaling keypoints to work reliably on richer datasets and environments is an important future area of research. Further, tracking objects over long temporal sequences can enable learning object dynamics and affordances which could be used to inform learning policies. A limitation of our model is that we do not currently handle moving backgrounds. Recent work [9] that explicitly reasons about camera / ego motion could be integrated to globally transport features between source and target frames. In summary, our experiments provide clear evidence that it is possible to learn visual abstractions and use simple algorithms to produce highly data efficient control policies and exploration procedures.

Acknowledgements. We thank Loic Matthey and Relja Arandjelović for valuable discussions and comments.

References

- [1] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 2013.
- [2] J. Buckman, A. Roy, C. Raffel, and I. Goodfellow. Thermometer encoding: One hot way to resist adversarial examples. 2018.
- [3] C. P. Burgess, L. Matthey, N. Watters, R. Kabra, I. Higgins, M. Botvinick, and A. Lerchner. Monet: Unsupervised scene decomposition and representation. *arXiv preprint arXiv:1901.11390*, 2019.
- [4] X. Chen, Y. Duan, R. Houthoofd, J. Schulman, I. Sutskever, and P. Abbeel. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2016.
- [5] A. Ecoffet, J. Huizinga, J. Lehman, K. O. Stanley, and J. Clune. Go-explore: a new approach for hard-exploration problems. *arXiv preprint arXiv:1901.10995*, 2019.
- [6] C. Finn, X. Y. Tan, Y. Duan, T. Darrell, S. Levine, and P. Abbeel. Deep spatial autoencoders for visuomotor learning. In *Proceedings of the International Conference on Robotics and Automation (ICRA)*, 2016.
- [7] V. Goel, J. Weng, and P. Poupart. Unsupervised video object segmentation for deep reinforcement learning. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [8] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2014.
- [9] A. Gordon, H. Li, R. Jonschkowski, and A. Angelova. Depth from videos in the wild: Unsupervised monocular depth learning from unknown cameras. *arXiv preprint arXiv:1904.04998*, 2019.
- [10] K. Greff, R. L. Kaufman, R. Kabra, N. Watters, C. Burgess, D. Zoran, L. Matthey, M. Botvinick, and A. Lerchner. Multi-object representation learning with iterative variational inference. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2019.
- [11] M. Grundmann, V. Kwatra, M. Han, and I. Essa. Efficient hierarchical graph-based video segmentation. In *Proceedings of the conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- [12] M. Hessel, J. Modayil, H. Van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot, M. Azar, and D. Silver. Rainbow: Combining improvements in deep reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- [13] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.

- [14] D. Horgan, J. Quan, D. Budden, G. Barth-Maron, M. Hessel, H. Van Hasselt, and D. Silver. Distributed prioritized experience replay. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- [15] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *CoRR*, 2015. URL <http://arxiv.org/abs/1502.03167>.
- [16] C. Ionescu, T. Kulkarni, A. van den Oord, A. Mnih, and V. Mnih. Learning to Control Visual Abstractions for Structured Exploration in Deep Reinforcement Learning. In *NeurIPS Deep Reinforcement Learning Workshop*, 2018.
- [17] T. Jakab, A. Gupta, H. Bilen, and A. Vedaldi. Unsupervised learning of object landmarks through conditional image generation. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [18] X. Ji, J. F. Henriques, and A. Vedaldi. Invariant information distillation for unsupervised image segmentation and clustering. *arXiv preprint arXiv:1807.06653*, 2018.
- [19] L. Kaiser, M. Babaeizadeh, P. Milos, B. Osinski, R. H. Campbell, K. Czechowski, D. Erhan, C. Finn, P. Kozakowski, S. Levine, et al. Model-based reinforcement learning for atari. *arXiv preprint arXiv:1903.00374*, 2019.
- [20] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [21] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2012.
- [22] T. D. Kulkarni, W. F. Whitney, P. Kohli, and J. Tenenbaum. Deep convolutional inverse graphics network. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2015.
- [23] T. D. Kulkarni, K. Narasimhan, A. Saeedi, and J. Tenenbaum. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2016.
- [24] S. Li, B. Seybold, A. Vorobyov, A. Fathi, Q. Huang, and C.-C. Jay Kuo. Instance embedding transfer to unsupervised video object segmentation. In *Proceedings of the conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [25] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [26] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 2015.
- [27] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell. Curiosity-driven exploration by self-supervised prediction. In *ICML*, 2017.
- [28] P. O. Pinheiro, T.-Y. Lin, R. Collobert, and P. Dollár. Learning to refine object segments. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [29] M. Plappert, R. Houthoofd, P. Dhariwal, S. Sidor, R. Y. Chen, X. Chen, T. Asfour, P. Abbeel, and M. Andrychowicz. Parameter space noise for exploration. *arXiv preprint arXiv:1706.01905*, 2017.
- [30] M. Riedmiller. Neural fitted q iteration—first experiences with a data efficient neural reinforcement learning method. In *European Conference on Machine Learning*. Springer, 2005.
- [31] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

- [32] Z. Shu, M. Sahasrabudhe, R. Alp Guler, D. Samaras, N. Paragios, and I. Kokkinos. Deforming autoencoders: Unsupervised disentangling of shape and appearance. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [33] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 2016.
- [34] E. S. Spelke and K. D. Kinzler. Core knowledge. *Developmental science*, 2007.
- [35] R. S. Sutton, D. Precup, and S. Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 1999.
- [36] S. Suwajanakorn, N. Snavely, J. J. Tompson, and M. Norouzi. Discovery of latent 3d keypoints via end-to-end geometric reasoning. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- [37] Y. Tassa, Y. Doron, A. Muldal, T. Erez, Y. Li, D. d. L. Casas, D. Budden, A. Abdolmaleki, J. Merel, A. Lefrancq, et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.
- [38] J. Thewlis, H. Bilen, and A. Vedaldi. Unsupervised learning of object landmarks by factorized spatial embeddings. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2017.
- [39] W. F. Whitney, M. Chang, T. Kulkarni, and J. B. Tenenbaum. Understanding visual concepts with continuation learning. *arXiv preprint arXiv:1602.06822*, 2016.
- [40] O. Wiles, A. Koepke, and A. Zisserman. Self-supervised learning of a facial attribute embedding from video. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [41] T. Xue, J. Wu, K. Bouman, and B. Freeman. Visual dynamics: Probabilistic future frame synthesis via cross convolutional networks. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2016.
- [42] Y. Zhang, Y. Guo, Y. Jin, Y. Luo, Z. He, and H. Lee. Unsupervised discovery of object landmarks as structural representations. In *Proceedings of the conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

Appendix A Implementation Details

The feature extractor Φ is a convolutional neural network with six Conv-BatchNorm-ReLu layers [15]. The filter size for the first layer is 7×7 with 32 filters, and 3×3 for the rest with number of filters doubled after every two layers. The stride was set to 2 for layer 3 and 5 (1 for the rest). *KeyNet* Ψ has a similar architecture but includes a final 1×1 regressor to K feature-maps corresponding to K keypoints. 2D coordinates are extracted from these K maps as described in [17]. The architecture of *RefineNet* is the transpose of Φ with $2 \times$ bilinear-upsampling to undo striding. We specify K for each environment but keep all other hyper-parameters of the network fixed across experiments. We used the Adam optimizer [20] with a learning rate of 0.001 (decayed by 0.95 every 10^5 steps) and batch size of 64 across all experiments.

For evaluating keypoint predictions, the distance threshold value (ϵ) was set to the average ground-truth spatial extent of entities for each environment given in the table below. Note, these ϵ values correspond to coordinates normalised to the $[-1, 1]$ range for both the (x, y) dimensions.

environment	pong	ms_pacman	frostbite	stack_4
ϵ	0.20	0.15	0.20	0.15

Code for the *Transporter* model is available at:

<https://github.com/deepmind/deepmind-research/tree/master/transporter>.

Appendix B Diverse Data Generation

To train the Transporter, a dataset of observation pairs is constructed from environment trajectories. It is important that this dataset contains a diverse range of situations, and unconditionally storing a pair from all trajectories generated by a random policy may contain many similar pairs. To mitigate this, we use a *diverse* data generation procedure as follows.

We generate trajectories of up to length 100 (action repeat is set to 4, so these trajectories represent up to 400 environment frames) using a uniform random policy, and uniformly sample one observation from the first half of the trajectory and one frame from the second half. Trajectories are shorter than 100 only when the end of an episode is reached. A buffer containing the maximum number of pairs we want to generate (in these experiments, 100k) is populated unconditionally from a number of environment actor threads until it is full. More frame pairs are generated, up to some defined maximum budget, and are conditionally written into the buffer as follows.

First some number of indices of existing pairs are sampled from the buffer, and for each of them we compute the nearest neighbor by L2 distance to other elements of the buffer. We take the same number of new generated frame pairs, and also compute their nearest neighbor in the buffer. For corresponding pairs of (existing frame pair, new frame pair) we overwrite the existing pair with the new pair whenever the new pair has a greater nearest neighbor distance, or if a uniform random number $\in [0, 1] < 0.05$. We continue this procedure until the maximum budget is reached, and write out the final buffer as our training set. Note that the reward function is not used at all in this procedure.

For efficiency, we store a lower resolution copy of the buffer (64×64 , grayscale) on the GPU to perform efficient nearest neighbor calculations, keeping corresponding higher res (128×128 RGB) copies on CPU RAM. Using a single consumer GPU and a 56 core desktop machine, with many actor subprocesses, this approach can perform 10 million environment steps (40 million total frames) in approximately 1 hour.

Appendix C Videos

Videos visualising various aspects of the model are available at:

https://www.youtube.com/playlist?list=PL3LT3tVQRpbvGt5fgp_bKGvW23jF11Vi2

Appendix D Pixel Transport versus Feature Transport

We investigated whether learning features is as important as spatially transporting them between frames. As shown in fig. 8, we show that transporting learned features significantly outperforms transporting pixels. Transporting pixels gives rise to ambiguous intermediate pixel representations, making it difficult for the final CNN decoder network to solve the downstream pixel prediction task. On the other hand, the feature encode higher level information and the decoder network learns a more abstract function to solve the prediction problem.

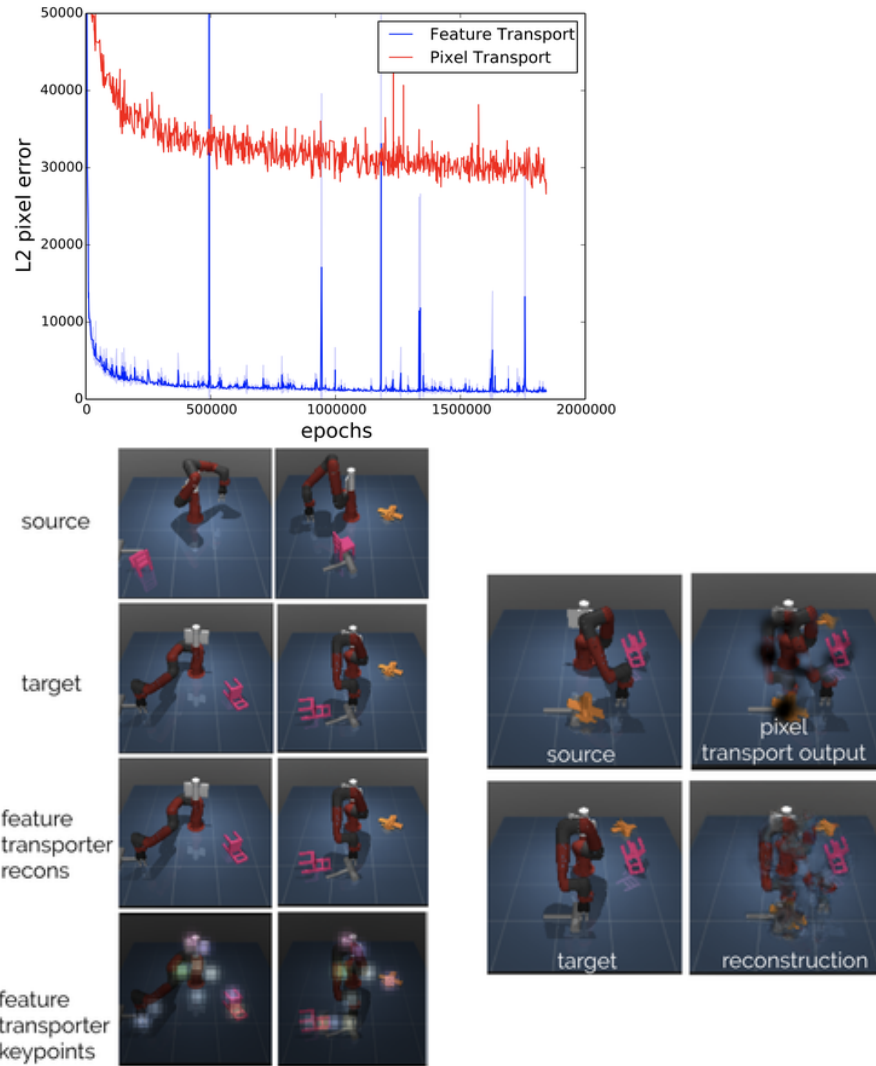


Figure 8: **Transporting features is significantly better than transporting pixels.** Given a sawyer arm with tabletop toys, *Transporter* discovers keypoints at the joint locations of the robot and object centroids (left two columns). In this experiment we investigate whether it is important to transport learned features or pixels. In case of pixel transport, the refinement network has to perform difficult and ambiguous computations to predict the target frame. Therefore the final pixel reconstruction error is significantly higher for pixel transporter (right most column).

Appendix E Temporal consistency of keypoints

Figures 9, 10, 11, and 12 show the inferred keypoints on frames selected from a single episode on Atari ALE [1] (Pong, Frostbite and Ms. Pac-Man) and Manipulator [37] (stack_4) domains. The selected frames are each 10 time steps apart. The first frame has been explicitly chosen to ensure there is enough diversity in the shown frames. The colours are time consistent – a specific colour corresponds to the same keypoint throughout the episode. Thus, if a given game ‘object’ is assigned the same coloured keypoint throughout the episode, that keypoint is temporally consistent for that ‘object’.

Videos showing the inferred keypoints by the three methods for entire episodes can be accessed at: https://www.youtube.com/playlist?list=PL3LT3tVQRpbvGt5fgp_bKGvW23jF11Vi2

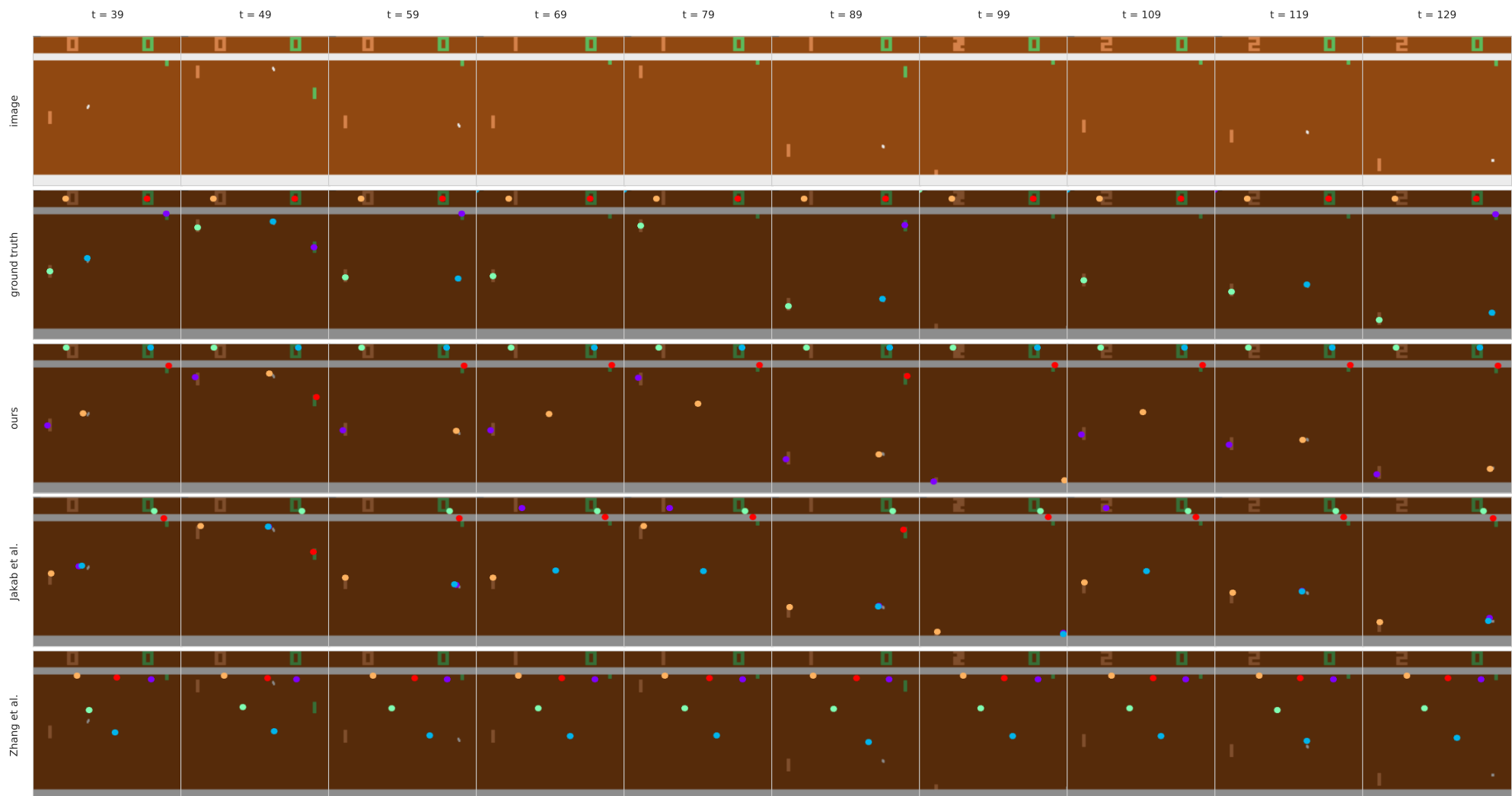


Figure 9: Atari [1]: pong



Figure 10: Atari [1]: frostbite

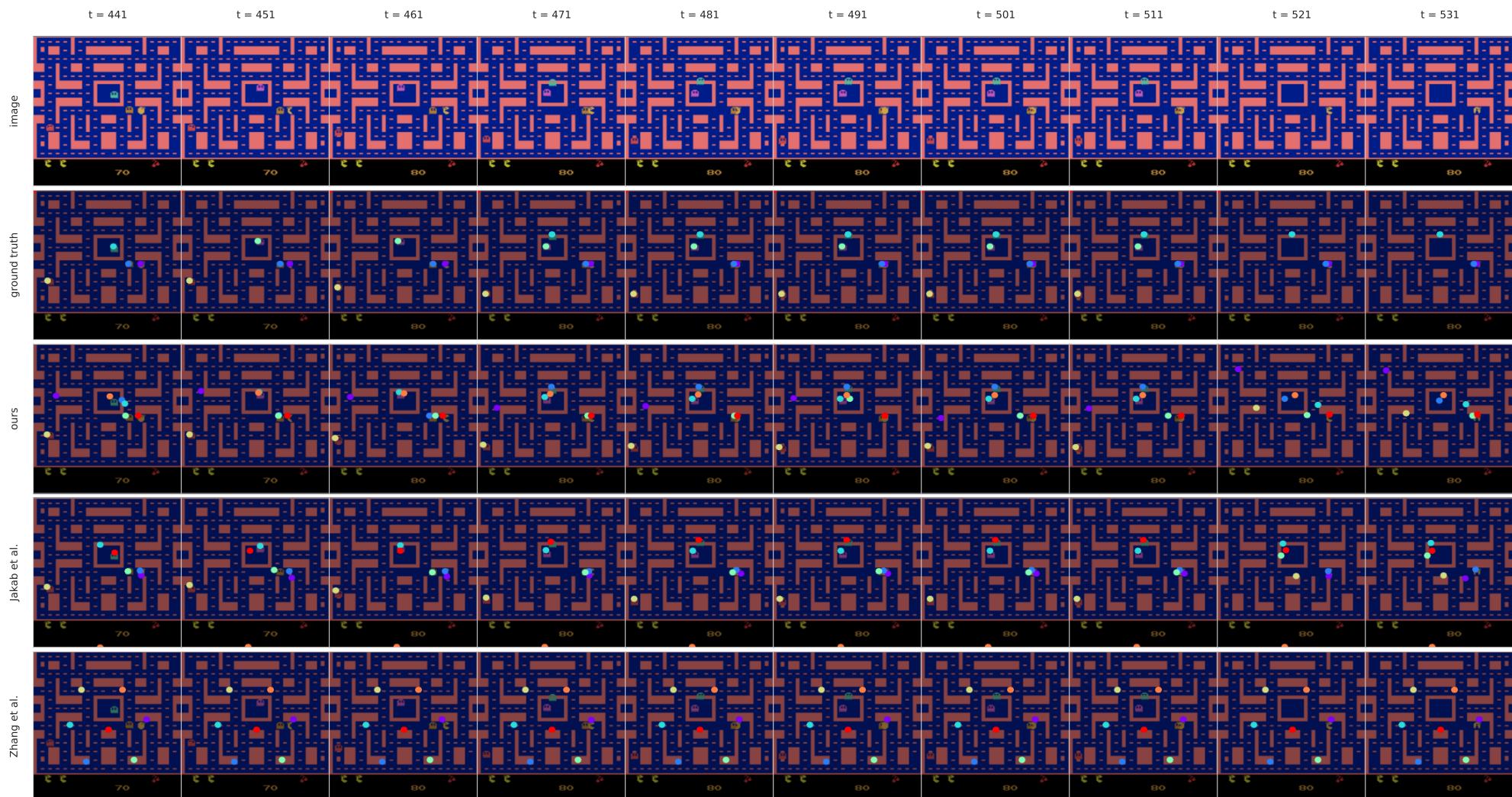


Figure 11: Atari [1]: ms pacman

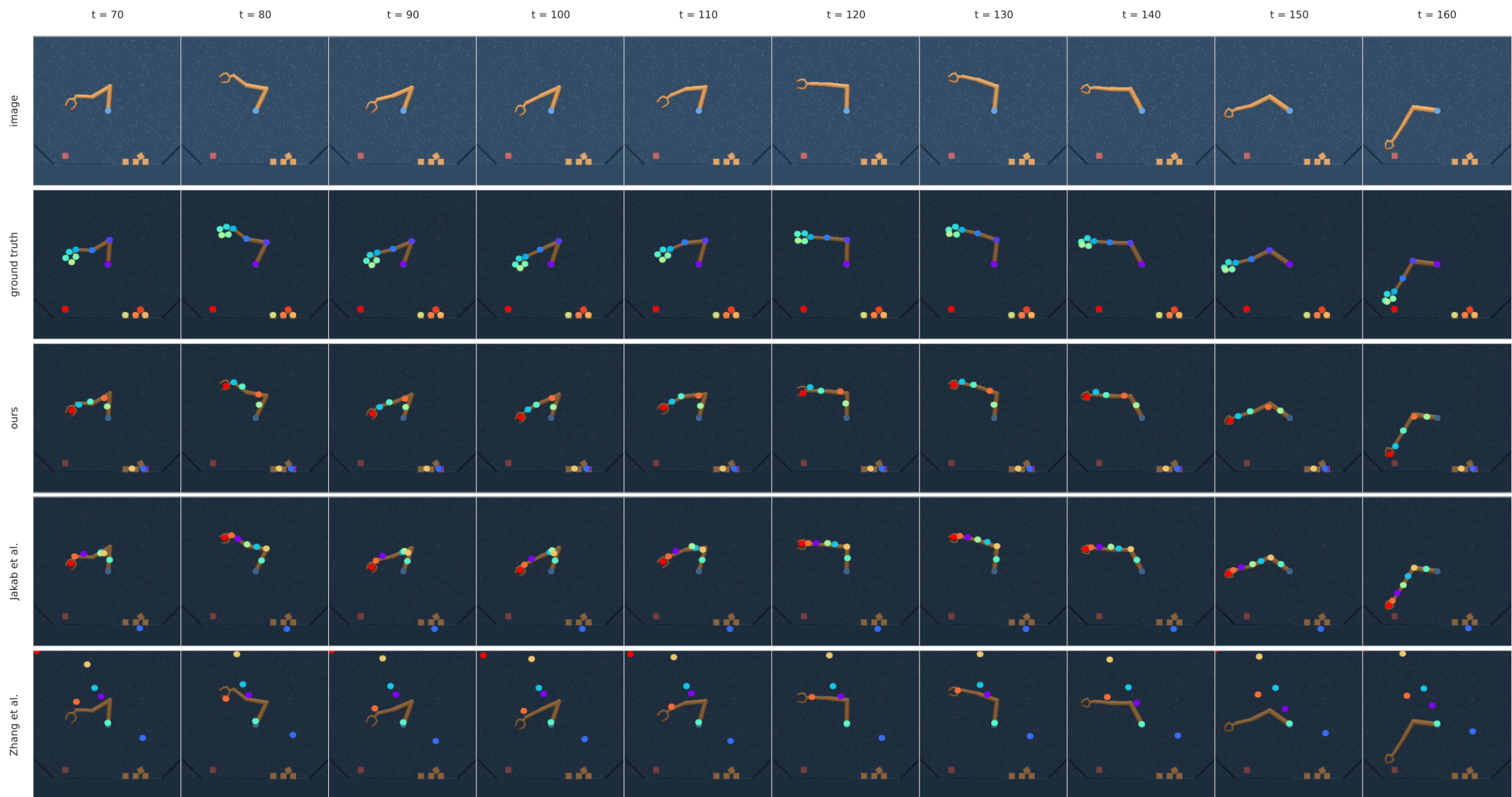


Figure 12: Manipulator [37]: stacker with 4 objects

Appendix F Reconstructions

We visualise the reconstructed images on Atari ALE [1] (Figures 13, 14, and 15) and Manipulator [37] domains (Figure 16) for randomly selected frames. The rows in the figures correspond respectively to our model, Jakab and Gupta et al. [17] and Zhang et al. [42]. The first two columns are the inputs given to the models. Whereas our model requires a pair of input frames (*image* and *future_image*), the remaining two models only require single frame (*future_image*). The third column (*reconstruction*) shows the reconstructed target image. The final column (*keypoints*) shows the inferred keypoints for the given inputs.

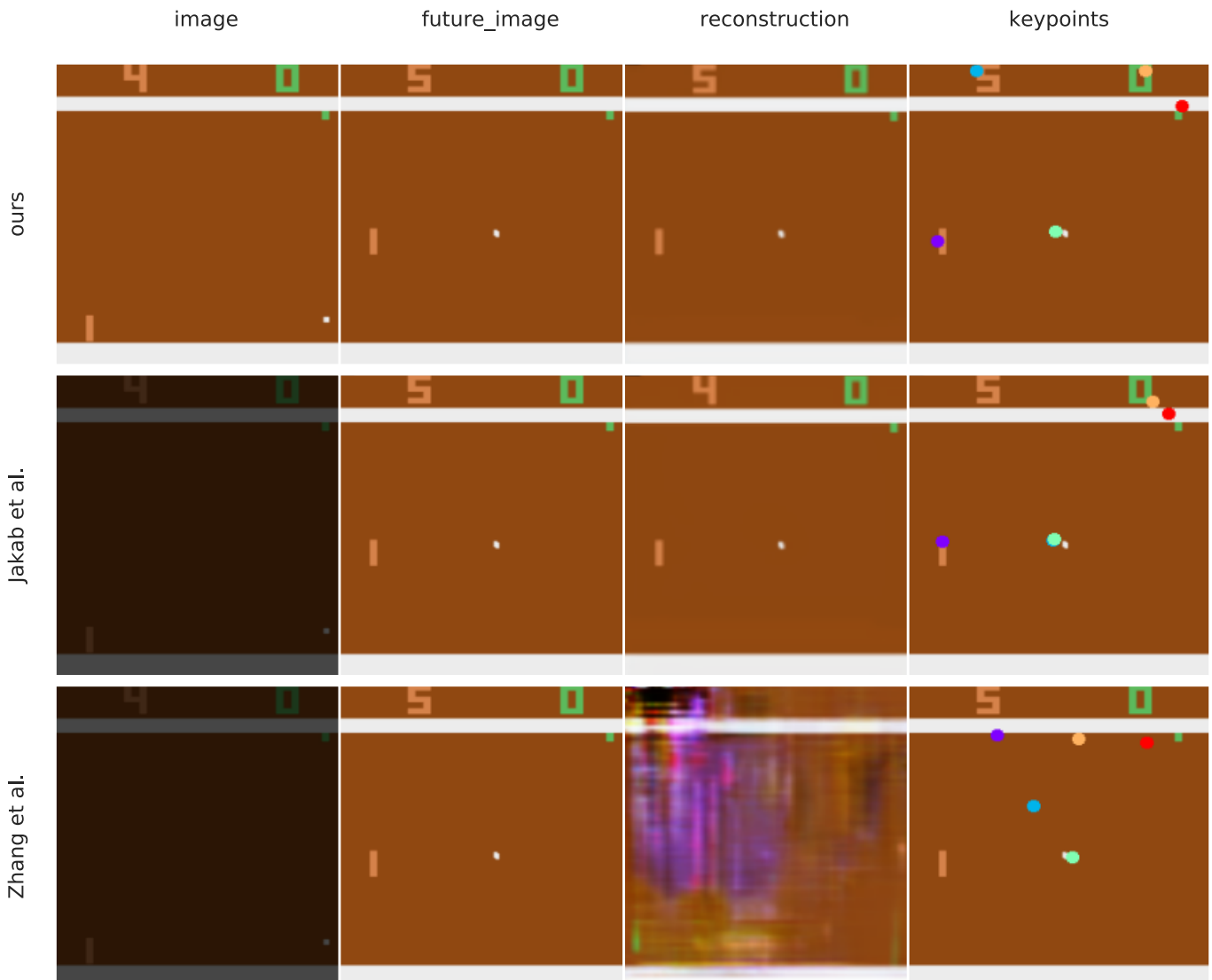


Figure 13: Reconstruction: pong

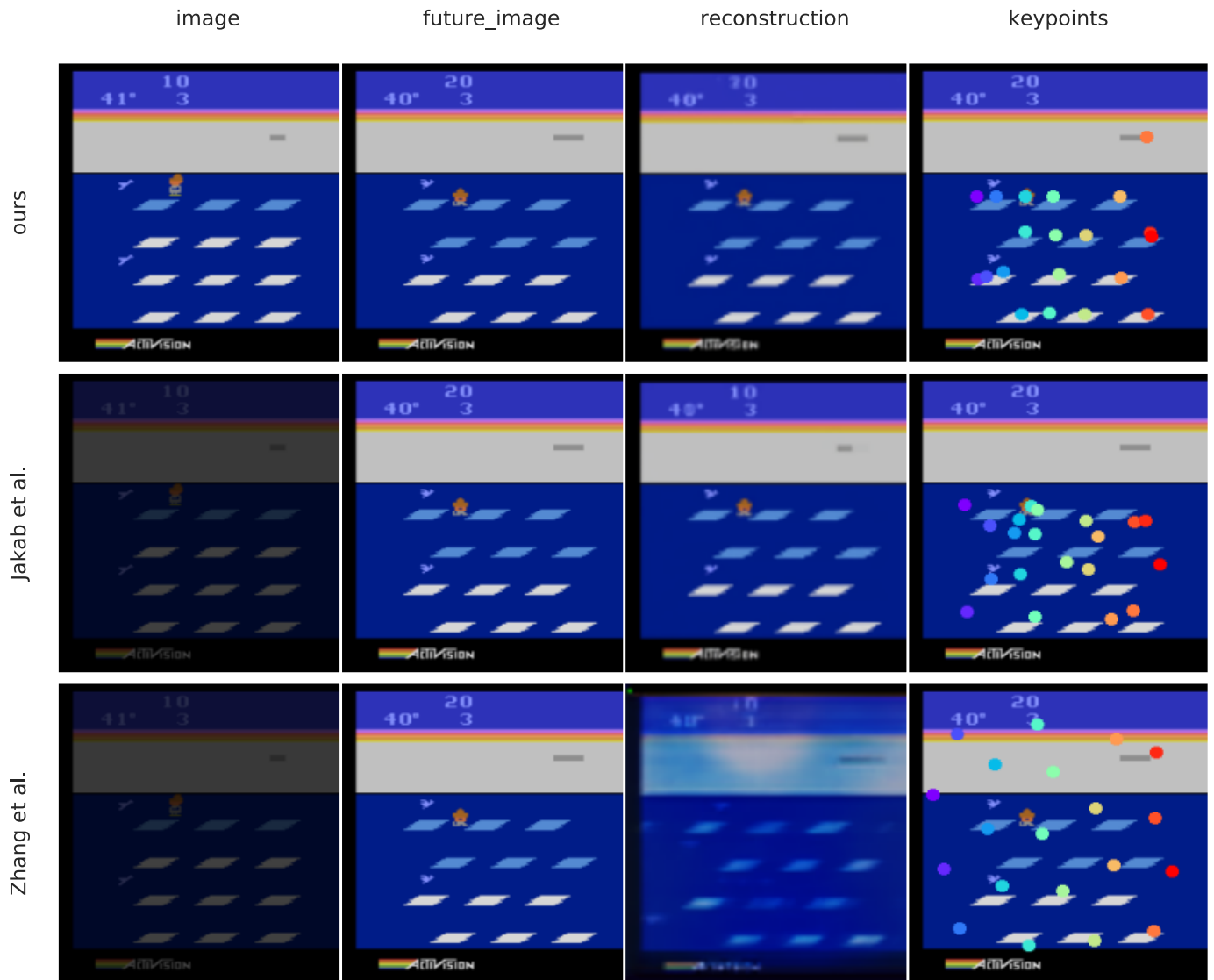


Figure 14: Reconstruction: frostbite



Figure 15: Reconstruction: ms pacman

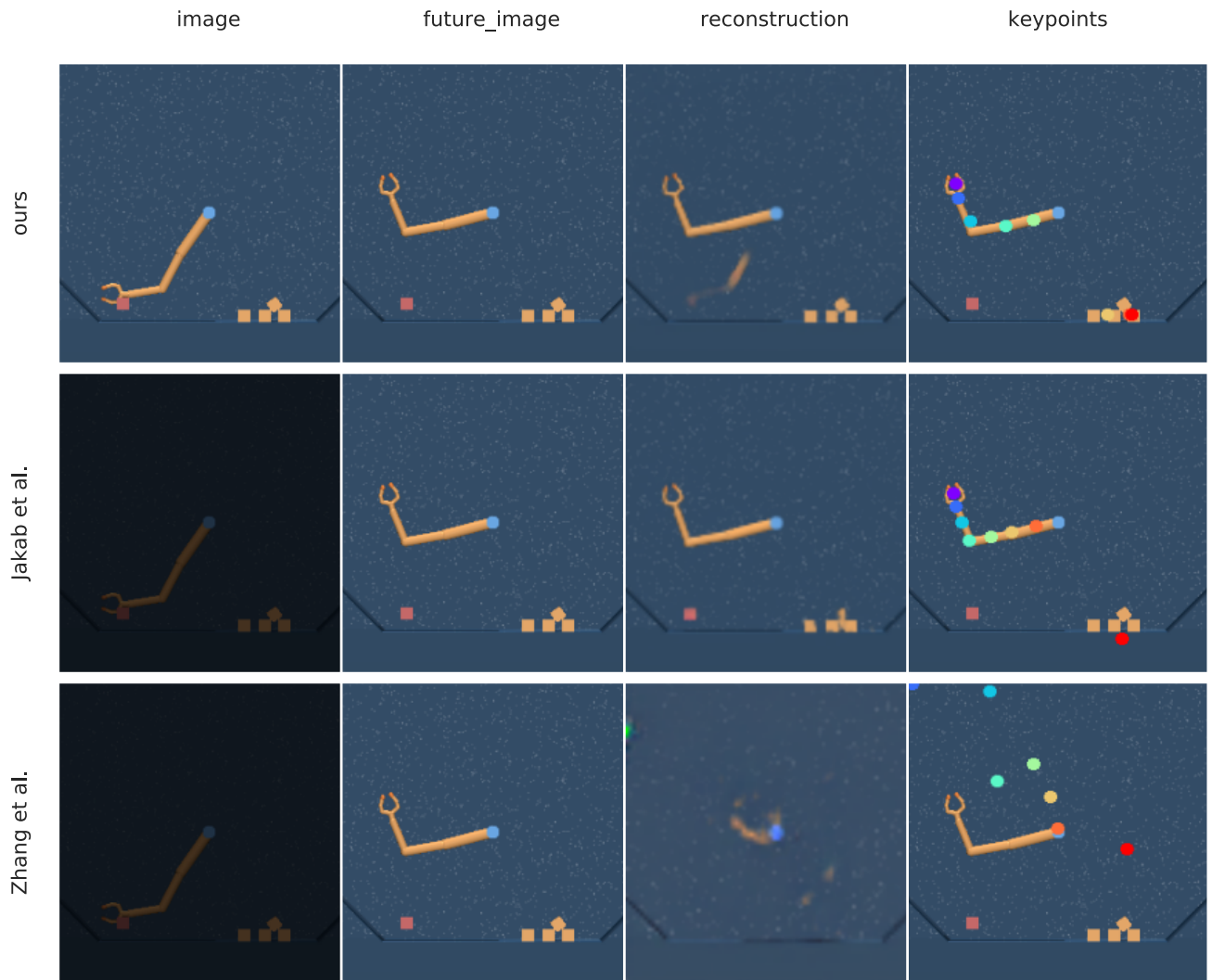


Figure 16: Reconstruction: stacker with 4 objects