

DISCRETIZATION OF CONTINUOUS INPUT SPACES IN THE HIPPOCAMPAL AUTOENCODER

Anonymous authors

Paper under double-blind review

ABSTRACT

Understanding the encoding mechanisms of hippocampal place cells remains a significant challenge in neuroscience. Although sparse autoencoders have been shown to exhibit place cell-like activity, the underlying processes are not fully understood. In this study, we compare spatial representations learned by dense and sparse autoencoders trained on images of 3D environments and find that only sparse autoencoders with orthonormal activity regularization in latent space produce place cells. We then show that this regularization promotes similar images to map onto the same neurons, acting as a locality-sensitive hash function. Notably, we demonstrate that these neurons are visually interpretable through activity clamping and decoding, suggesting the formation of detailed episodic memories at the single-neuron level. We then introduce a novel metric to quantify how neurons discretize the image space into disjoint receptive fields, revealing that sparse autoencoders tile input spaces with minimal overlap. Furthermore, we observe that whereas dense autoencoders generate population codes resembling visual cortex activity near criticality, sparse autoencoders produce higher-dimensional codes, thus suggesting a similar coding strategy in the hippocampus. Extending our approach to the auditory domain, we also replicate the emergence of "frequency place cells" by training sparse autoencoders on audio snippets sampled from a frequency-varying signal, and show that population representations retain the statistical structure of the sample distribution. Lastly, we demonstrate that reinforcement learning agents can leverage these high-dimensional image representations to solve complex spatial-cognitive tasks, despite their inherent brittleness. Overall, our findings elucidate how sparse input compression in autoencoders can give rise to discrete, interpretable memories, establishing an explicit link between episodic memory formation and spatial representations in the hippocampus.

1 INTRODUCTION

Early physiological experiments with rats revealed that certain neurons in the hippocampus exhibit increased activity when the animal occupies specific regions of the environment (O'Keefe & Dostrovsky (1971)). Ever since the discovery of such "place cells", decades of animal research have established the hippocampus as a neural system that learns a cognitive map of the environment and uses it for spatial navigation (Moser et al., 2008). Subsequent experimental studies also identified the hippocampus as a key structure in episodic memory formation (Moser et al., 2015). Although several attempts have been made to unify these observations under a coherent conceptual framework (Redish, 1999; Eichenbaum, 2017), a clear mechanistically relationship episodic memory and spatial representations remains elusive. Furthermore, numerous experiments reporting the instability of place cell activity over time and their modulation by other non-spatial variables (Fenton & Muller, 1998; Jercog et al., 2019) raise an open question: what are these cells truly encoding?

Efforts to answer this question have demonstrated that place cell-like activity can emerge under various conditions: when artificial agents optimize a predictive coding objective (Recanatesi et al., 2021; Uria et al., 2020; Ratzon et al., 2023; Gornet & Thomson, 2023; Levenstein et al., 2024; Chen et al., 2022), when networks optimize temporal stability and pairwise decorrelation in processing visual inputs (Wyss et al., 2006), or when building sparse, compressed representations of environmental states (Santos-Pata et al., 2021a;b; Benna & Fusi, 2021; Ketz et al., 2013). Notably, the approach where sparse compression of information leads to spatial tuning aligns with the earlier hippocampal

054 autoencoder model (Gluck & Myers, 1993), and has been shown to replicate several distinct place
 055 cell phenomena following environmental manipulations (Santos-Pata et al., 2021a;b).

056
 057 In this work, we further investigate the mechanisms behind episodic memory formation and the
 058 emergence of place cells in the hippocampal autoencoder model. We demonstrate that sparse
 059 autoencoders equipped with orthonormal activity regularization can create discontinuities in the
 060 manifold of the latent space, discretizing arbitrary input spaces into disjoint receptive fields, whereby
 061 subsets of similar inputs converge onto distinct neurons. When applied to visual images, this clustering
 062 process generates place cells operating on a very high-dimensional population code. In turn, these
 063 neurons are shown to encode detailed visual memories. Moreover, we show that similar effects result
 064 from applying the same principle in the auditory domain, recapitulating recently reported "frequency
 065 place cells". Lastly, we show that reinforcement learning agents can make use of such sparse and
 066 high-dimensional hippocampal-like representations to solve spatial-cognitive tasks.

067 2 MODEL AND RESULTS

068 2.1 HIPPOCAMPAL-LIKE PLACE CELLS EMERGE IN SPARSE AUTOENCODERS

069 We studied the learning of spatial representations by training autoencoders (Figure 1a) with randomly
 070 sampled images from four different tasks in the Animal-AI environment (Beyret et al., 2019): Double
 071 T-maze, Cylinder, Object Permanence, and Thorndike. We trained two types of autoencoders. "Dense"
 072 autoencoders aimed solely to reconstruct the input images, thus preserving input information in
 073 their latent space Z . In contrast, "sparse" autoencoders had an additional objective beyond input
 074 reconstruction: to develop sparse activity patterns in the latent space Z . This was achieved using the
 075 following loss function:
 076

$$077 \mathcal{L} = \frac{1}{m} \|\mathbf{X} - \hat{\mathbf{X}}\|_2^2 + \frac{\lambda}{mn} \|\mathbf{I}_n - \mathbf{Z}^T \mathbf{Z}\|_F, \quad (1)$$

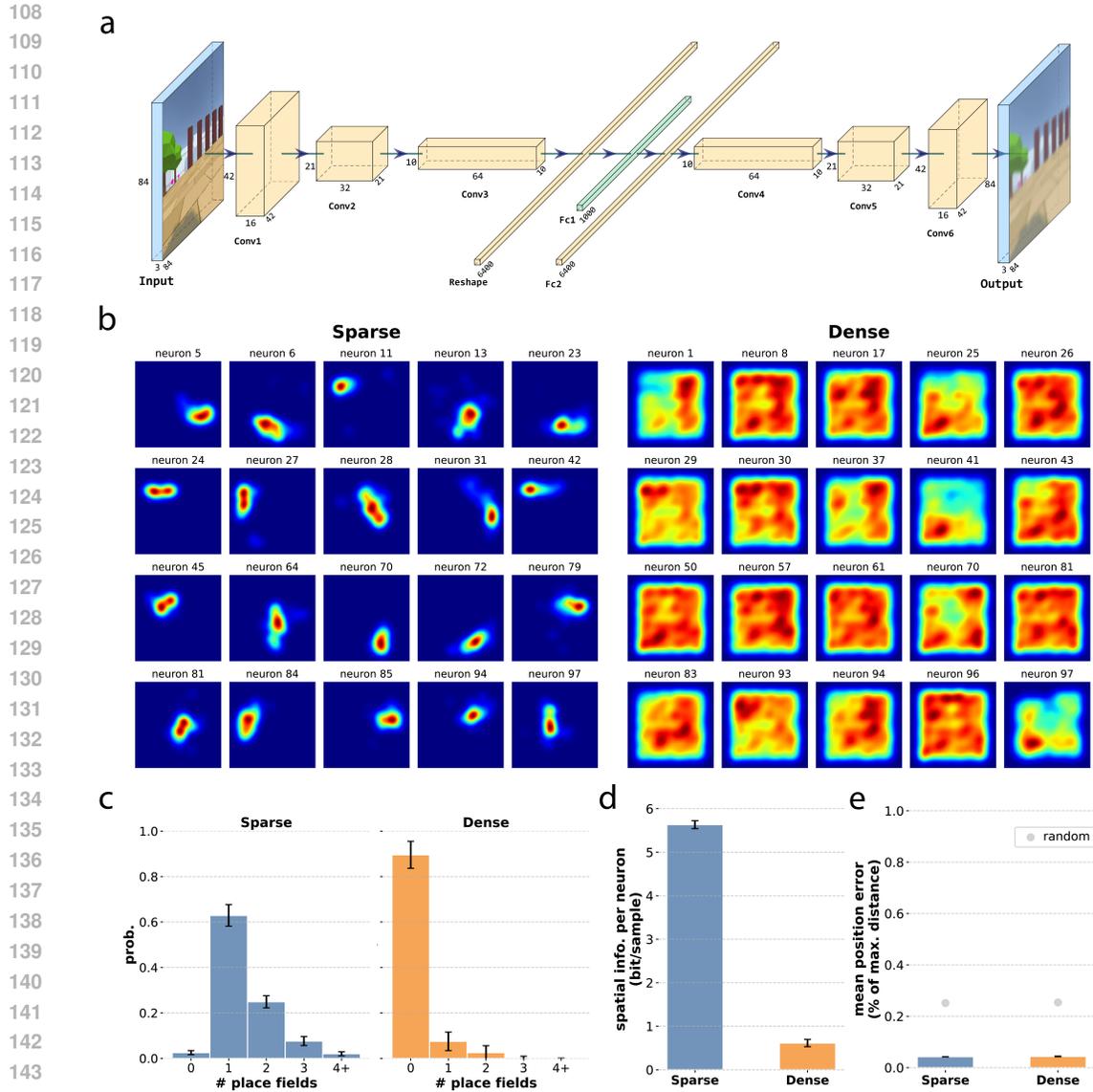
078 where m denotes the batch size, n the number of neurons in Z , and the first term is the mean
 079 squared error (MSE) between inputs \mathbf{X} and their reconstructions $\hat{\mathbf{X}}$, encouraging Z to preserve input
 080 information. The second term is an orthonormal activity regularization term, whose strength is
 081 controlled by λ , pushing the Gramian $\mathbf{Z}^T \mathbf{Z}$ towards the identity matrix \mathbf{I}_n . Since $\mathbf{Z}^T \mathbf{Z}$ captures
 082 the co-activation strengths between neurons in a training sample batch, the orthonormal activity
 083 regularization promotes pairwise decorrelation while ensuring equal contribution across neurons.
 084 We found this approach yields improved and more reliable results compared to the L1 activity
 085 regularization term typically used in sparse autoencoders, particularly in alleviating the dead ReLU
 086 problem (Lu et al., 2019). For dense autoencoders, λ was set to zero, leaving only the reconstruction
 087 error. We refer the reader to the Detailed methods section in the Appendix for a complete description
 088 of the environments, dataset generation, and parameters used in this study.

089
 090 Training both types of autoencoders yielded significantly different internal representations of space
 091 in their latent space. Dense autoencoders developed many neurons that fired almost everywhere in
 092 space, with no defined place fields. In contrast, sparse autoencoders developed a majority of neurons
 093 with one or two localized place fields, similar to place cells in the hippocampus (Figure 1b, c). The
 094 spatial specificity of sparse autoencoder neurons was also reflected in significantly higher spatial
 095 information scores compared to dense autoencoder neurons (Figure 1d). These results demonstrate
 096 that single-unit spatial tuning emerges in sparse autoencoders but not in dense autoencoders, despite
 097 both types of networks containing the same amount of positional information at the population level,
 098 as shown by linear decoding analyses (Figure 1e).

099 2.2 SPARSE AUTOENCODERS DISCRETIZE AND TILE THE IMAGE SPACE WITH INTERPRETABLE 100 NEURONS

101
 102 Identifying neurons with spatial selectivity similar to hippocampal place cells allowed us to investigate
 103 what these neurons encode. Given that their spatial selectivity must arise from some form of visual
 104 selectivity, we explored whether they also exhibit localized receptive fields in image space.

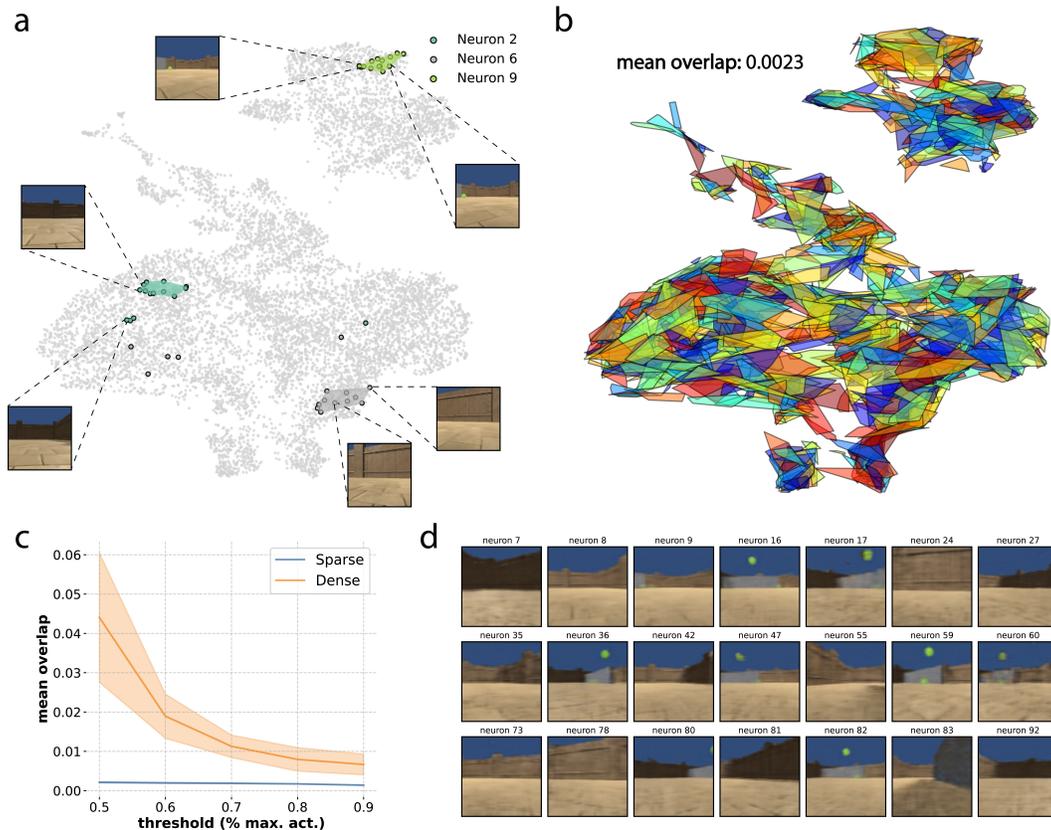
105
 106 We created an image space by extracting semantically-relevant image embeddings of all samples
 107 using CLIP and further reducing the dimensionality to a 2D space with UMAP. We then searched
 for clusters in this 2D image space by running the DBSCAN algorithm on the points corresponding



146 Figure 1: Hippocampal-like place cells emerge in sparse autoencoders. (a) Autoencoder architecture,
147 featuring the hidden layer or latent space Z (denoted as Fc1) with 1000 neurons. (b) Representative
148 examples of the neurons’ spatial ratemaps for sparse and dense autoencoders. (c) Probability
149 distribution of place field number across environments. (d) Average spatial information per neuron
150 across environments. (e) Normalized average distance error of linear decoding of position with
151 the ratemaps’ population vectors, across environments. The grey dots represent the expected linear
152 decoding errors after performing 1000 random permutations of the ratemaps’ values.

153
154
155
156 to images that maximally activated a particular neuron (see Figure 2a and Detailed methods in the
157 Appendix). These clusters formed convex hulls (i.e., patches) that corresponded to the neuron’s
158 receptive fields in the image space. When pooled together, receptive fields across neurons partitioned
159 and covered the entire image space (see example in Figure 2b). Furthermore, we computed an overlap
160 metric to estimate the redundancy across the neurons’ receptive fields. We observed that sparse
161 autoencoder neurons tiled the image space in a minimally-overlapping manner, in contrast to dense
autoencoder neurons, whose overlap tended to be significantly higher (Figure 2c).

162 Additionally, we performed unit clamping experiments, setting neurons to their maximal recorded
 163 value while others were set to zero, and then decoded their activity back to images. The generated im-
 164 ages showed a striking resemblance to the training images, making these neurons highly interpretable
 165 (Figure 2d). These results establish a solid relationship between episodic memory formation and
 166 spatial coding.



196 Figure 2: Sparse autoencoders discretize and tile the image space with interpretable neurons. (a)
 197 Images taken in one of the environments ('Cylinder'), encoded with CLIP and further reduced to two
 198 dimensions with UMAP. Points of different colors correspond to the images that maximally activate
 199 each example neuron (above the 50% threshold of the maximum neuron's recorded activity). Clusters
 200 of maximally activated images are extracted with DBSCAN, making up the convex hulls. (b) Convex
 201 hulls for all neurons in a sparse autoencoder trained with images from the 'Cylinder' environment.
 202 The overlap metric corresponds to the expected overlapping area (in %) of two randomly chosen
 203 hulls (see Detailed methods in the Appendix for further details). (c) Average overlap in 2D image
 204 space of sparse and dense autoencoders, across tasks and for a range of threshold values of maximal
 205 activation. (d) Example interpretable neurons in the sparse autoencoder. The corresponding neuron
 206 in latent space Z is set to its maximum recorded value across the dataset, while all other neurons are
 207 set to zero. Then, the enforced activity vector Z is deconvolved into an image by passing it through
 208 the decoder.

209 2.3 HIGH-DIMENSIONAL POPULATION STRUCTURE IN SPARSE AUTOENCODERS

211 Having linked the formation of episodic memories with the discretization of the image space in sparse
 212 autoencoders, we explored the population structure of the latent space representations. Inspired by
 213 Stringer et al. (2019) on the dimensionality of the population code in the mouse visual cortex, we
 214 examined the dimensionality in our autoencoders. Dimensionality was estimated by performing
 215 PCA on Z and computing the linear fit of the resulting eigenspectrum in log-log space, yielding a
 power-law exponent, α . High α values indicate low-dimensional codes, while low α values suggest

high-dimensional codes. An $\alpha \approx 1$ indicates a criticality regime with a high-dimensional but smooth (i.e., no discontinuities) underlying manifold, as seen in neural responses in the visual cortex (Stringer et al., 2019).

We found that dense autoencoders had dimensionality scores close to 1, similar to visual cortex (Stringer et al., 2019), whereas sparse autoencoders exhibited higher-dimensional representations (Figure 3a), aligning more with the efficient coding hypothesis (Barlow et al., 1961). The almost-flat eigenspectrum suggests that sparse autoencoders’ population activity indeed encodes fine stimulus features. Moreover, the orthonormal activity regularization also disrupted the input-output similarity preservation typically seen in dense autoencoders (Figure 3b), further supporting the idea of a sharp discretization of the image space by sparse autoencoders.

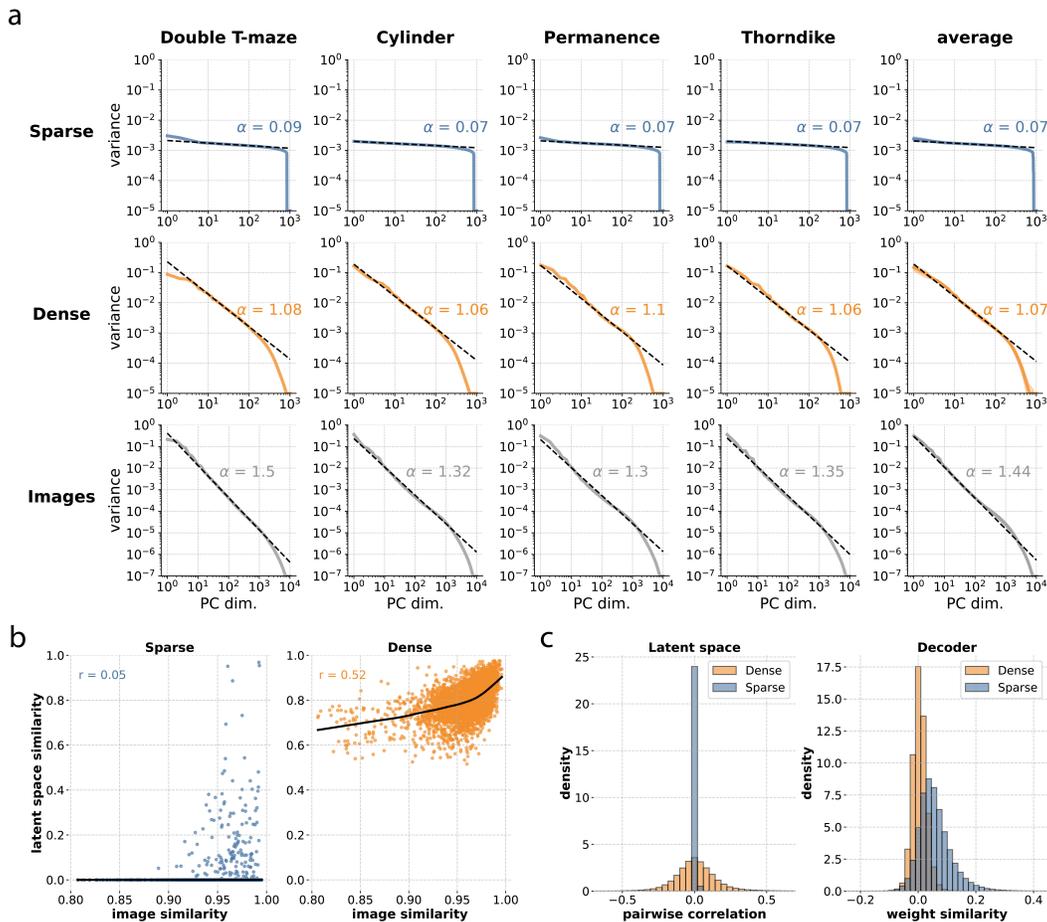


Figure 3: High-dimensional population structure in sparse autoencoders is grounded on mixed selectivity. (a) Eigenspectrum decay in latent space representations (first two rows) and images from the environments (third row). Parameter α corresponds to the power law exponent from linear fitting in log-log space. (b) Input-output similarity for sparse and dense autoencoders, with data pooled across environments. Correlation scores correspond to Spearman’s rank coefficients, and fitting curves have been generated with a locally weighted scatter-plot smoother (LOWESS) for improved visualization. (c) Pairwise Pearson correlation scores between all neurons’ activity in latent space, pooled across environments (left) and pairwise kernel similarity in the decoder weights (layer Fc2), representing the similarity density across ’words’ in the learned ’dictionary’ (right).

Borrowing concepts from sparse dictionary learning (Lewicki & Sejnowski, 2000), we considered the decoder weights to be the dictionary of kernels, and the sparse neuron activities to be the coefficients

that use the dictionary to reconstruct the inputs. Dense autoencoders exhibited orthogonal kernels, while sparse autoencoders showed non-orthogonal kernels despite highly uncorrelated activity patterns (Figure 3c). This indicates that orthogonal activity does not imply orthogonal kernels, and that sparse autoencoder neurons learned similar feature combinations, indicative of mixed selectivity in neurons (Fusi et al., 2016). These findings suggest that mixed feature selectivity and high dimensionality are closely linked to the formation of detailed episodic memories.

2.4 ZERO-SHOT LEARNING OF PLACE CELLS IN SPARSE AUTOENCODERS

A typical observation in the hippocampus is that place cells can be identified within the first minutes of an animal being exposed to a new environment (Frank et al., 2004). Given that we have shown that sparse autoencoders exhibit very high dimensionality (Figure 3a) and single neurons tend to encode small clusters of samples (Figure 2), we investigated the extent to which neurons that learned to encode samples in one environment could generalize to encoding unseen environments and exhibit zero-shot place cells. Therefore, we trained sparse autoencoders in one environment and tested them across all others. Strikingly, neurons developed place fields with distributions very similar to those in their training environment (Figure 4a). Furthermore, the average spatial information across neurons and the mean decoding error of the rate maps were very similar, with no significant degradation compared to the training environment (Figure 4b). These results suggest that the network’s circuitry learned to cluster similar samples onto single neurons in a more generic manner, beyond the specific details of the training data.

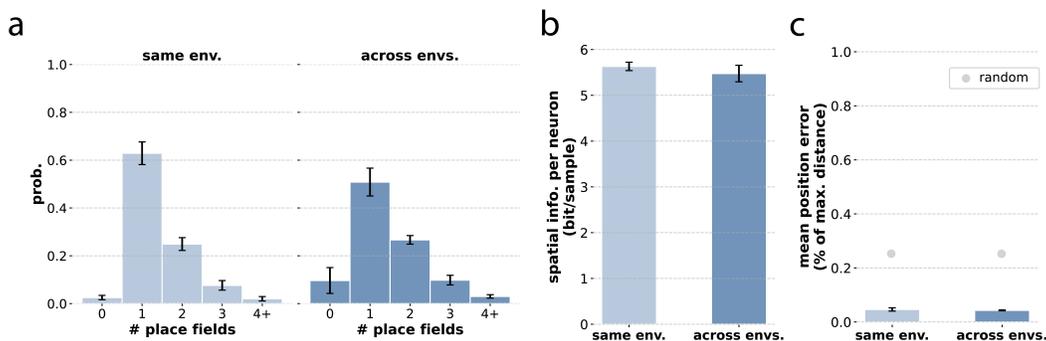


Figure 4: Zero-shot place cells in sparse autoencoders. (a) Probability distributions of place field number when testing a model within its training environment (light blue) or across unseen environments (dark blue). (b) Average spatial information per neuron, pooled across models and testing environments. (c) Normalized average distance error of linear decoding of position with the ratemaps’ population vectors, across models and testing environments. The grey dots represent the expected linear decoding errors after performing 1000 random permutations of the ratemaps’ values.

2.5 SPARSE AUTOENCODERS DISCRETIZE AND TILE THE INPUT FREQUENCY SPACE IN AN EXPERIENCE-DEPENDENT MANNER

If the hippocampus functions as a generic, modality-independent episodic memory system, our findings with the sparse autoencoder should generalize to other input modalities, such as sound. Indeed, “place cell”-like activity in the hippocampus has been reported for tasks involving “navigating” the sound frequency space, with neurons developing localized receptive fields around particular sound frequencies (Aronov et al., 2017). To investigate whether a similar effect could be observed within our framework, we trained autoencoders to compress and encode sound waves uniformly sampled from a linearly-varying frequency signal (Figure 5a).

We observed the emergence of frequency-specific receptive fields in sparse autoencoders, but not in dense autoencoders (Figure 5b), reproducing the main observations in Aronov et al. (2017). These receptive fields tiled the entire frequency space in a linear manner. However, when sampling was biased towards certain frequencies, the neurons’ receptive fields became denser and clustered around those frequencies, preserving the statistical structure of the sample distribution (Figure 5d).

Additionally, similar to our previous findings with visually interpretable neurons, we found that individual neurons encoded particular frequencies so that these representations were readily decodable via activation clamping (Figure 5c). These results demonstrate that sparse autoencoders can discretize arbitrary input spaces to support episodic memory formation.

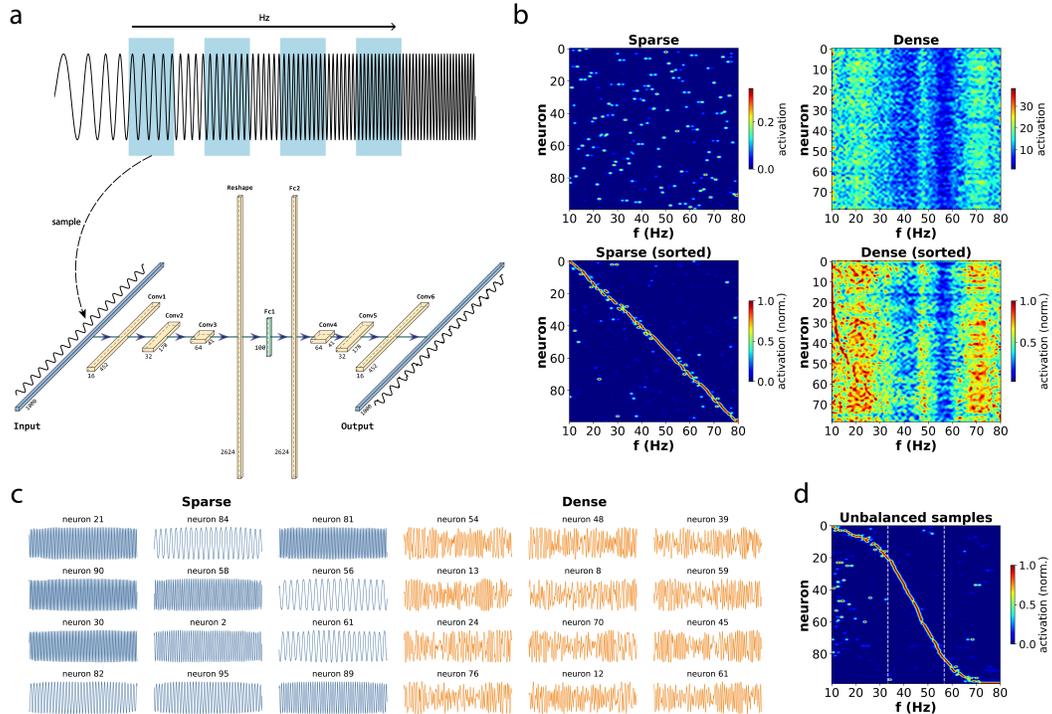


Figure 5: Sparse autoencoders discretize and tile the input frequency space in an experience-dependent manner. (a) Data samples are generated by applying a uniformly distributed sliding window to a linearly-varying frequency signal. The samples are fed into a convolutional autoencoder, analogous to the one used for vision (more details can be found in the Detailed Methods section of the Appendix). (b) Unsorted and sorted receptive fields by peak activity location for both sparse and dense autoencoders. Latent space activity Z responding to pure tone test inputs was convolved with a Gaussian kernel (sigma of 0.5 Hz), and then normalized by the maximum per neuron in the sorted plots. Lanczos interpolation was applied to all plots for improved visualization. (c) Decoded output signals after setting the corresponding neuron in latent space to its maximum recorded value across the dataset, while all other neurons were set to zero. (d) Sorted receptive fields in a sparse autoencoder trained with an unbalanced dataset. The data samples were generated with a sliding window that was not uniformly distributed in the frequency space, but whose density followed a Gaussian distribution centered at 45 Hz. Dashed vertical lines denote one standard deviation σ from the mean.

2.6 REINFORCEMENT LEARNING AGENTS LEARN EFFECTIVELY WITH SPARSE, HIGH-DIMENSIONAL REPRESENTATIONS

Representations used to build episodic memories are likely also employed for behavioral learning in the brain. Therefore, we investigated whether hippocampal-like representations emerging in sparse autoencoders would be suitable for reinforcement learning. To test this, we employed Deep Q-Networks (DQNs) (Mnih et al., 2015) incorporating either sparse or dense autoencoders to solve a range of tasks in the Animal-AI environment, which inherently require spatial navigation skills (see Figure 6a, and Detailed methods in the Appendix for further details on the tasks, model, and training parameters).

Very high-dimensional representations (such as those based on efficient coding) are known to be highly sensitive to slight input perturbations and are thought to generalize poorly to new, unseen

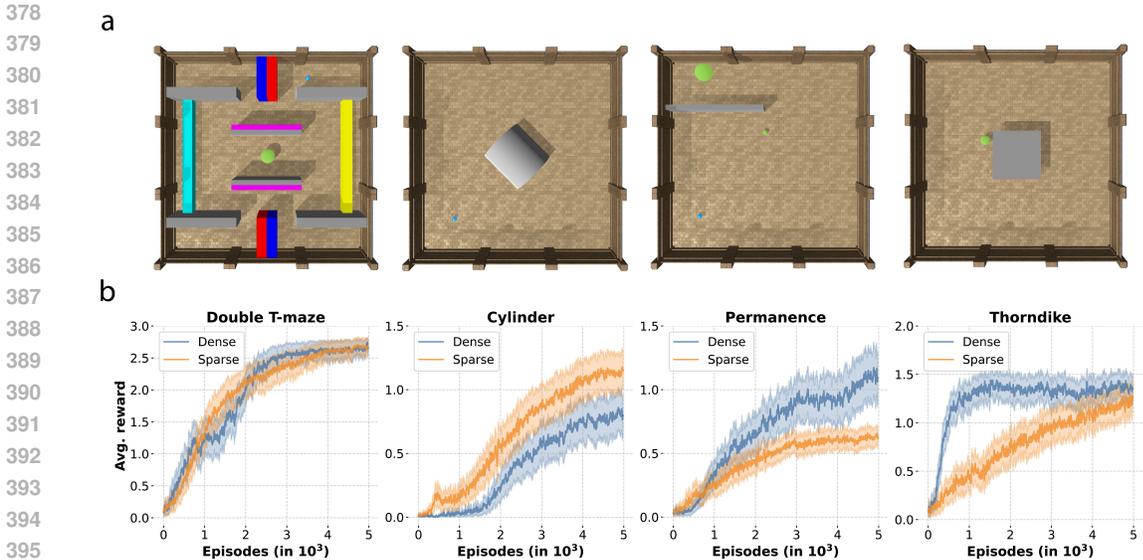


Figure 6: DQN agents learn effectively with sparse, very high-dimensional representations. (a) Overhead images captured from above the virtual arena of the four tasks (from left to right): Double T-maze, Cylinder, Object Permanence, and Thorndike. (b) Performance (average reward across episodes) of DQN agents using sparse and dense autoencoders across tasks.

samples (Nassar et al., 2020). Thus, one might expect that DQNs would struggle with tasks requiring generalization to unseen samples during training in non-stationary environments. Contrary to this expectation, we found that DQNs using sparse autoencoders were not systematically worse than those using dense autoencoders (Figure 6b). Although sparse autoencoders seemed to perform worse than dense autoencoders in two of the four tasks tested here (Object Permanence and Thorndike), they were superior in one of them (Cylinder) and had matched performance in the remaining one (Double T-maze). Therefore, while further testing is definitely needed to obtain a more reliable picture of their relative performance, these results suggest that, in practice, hippocampal-like representations can be suitable for reinforcement learning, despite their high dimensionality and inherent brittleness.

3 DISCUSSION

Optimization objectives underlying place cell emergence We have demonstrated the distinct emergence of place cells in autoencoders with orthonormal activity regularization (Figure 1). Notably, sparse compression alone was sufficient to develop spatial tuning, without the need for a predictive objective (Recanatesi et al., 2021; Ratzon et al., 2023; Uria et al., 2020; Gornet & Thomson, 2023; Levenstein et al., 2024; Chen et al., 2022). While predictive coding may explain other features of the hippocampus, such as place-cell theta sequences (Dragoi & Buzsáki, 2006), prediction does not appear to be necessary for the emergence of realistic place cells. We speculate that models optimized for next-input predictions likely learn compressed representations of the environment implicitly as part of the predictive objective. Furthermore, by training the autoencoders with randomly sampled and shuffled images, we have shown that neither temporal contiguity of samples nor any temporal correlations are required to develop place cells. This finding suggests that while predictive learning capturing temporal input correlations might correspond to experience-dependent theta sequences in the hippocampus, the formation of compressed state representations might correspond to time-independent learning processes at the gamma frequency scale (Lisman, 2005). Additionally, we have demonstrated that sparse autoencoders can learn localized receptive fields of the input space while breaking the relationship between input similarity and latent space similarity. This finding contrasts with previous research that emphasized the preservation of input-output similarity matching for learning localized receptive fields (Sengupta et al., 2018; Qin et al., 2023). Overall, it appears that

432 sparse compression alone is sufficient to learn localized receptive fields, that in turn manifest as place
433 cells when applied to the visual domain.
434

435 **Sparse and very high-dimensional codes** A slowly-decaying eigenspectrum, where fine details are
436 over-weighted, represents codes that create discontinuities by disrupting the locality of the manifold
437 structure supporting the input space distribution (Nassar et al., 2020). We demonstrated that this
438 discontinuity can be induced by an orthonormal activity regularization objective, facilitating the
439 creation of event memories by discretizing the image input space (Figure 2). Our results suggest
440 that very high-dimensional codes underlie the formation of place cells (Figure 3), aligning with the
441 efficient coding hypothesis (Barlow et al., 1961), which posits that the brain maximizes information
442 by eliminating correlations in sensory inputs, leading to sparse coding (Olshausen & Field, 1996).
443 Indeed, it has been shown that hippocampal neurons in rodents become sparser with prolonged
444 exposure to the environment (Ratzon et al., 2023). Moreover, the storage of social memories in
445 mice has been linked to high-dimensional representations in the hippocampus (Boyle et al., 2024).
446 Importantly, sparsity has been shown to control a generalization-discrimination trade-off (Barak
447 et al., 2013), which could explain why the hippocampus relies on sparse representations, well-suited
448 for progressive discrimination of similar environments and events. Our study shows that smooth
449 place cell maps can coexist with and emerge from extremely high-dimensional sparse codes (Chettih
450 et al., 2023). We therefore predict that the dimensionality of the population code along the sensory
451 hierarchy should decrease to support the learning of invariant sensory representations (Froudarakis
452 et al., 2020), and then increase sharply at the apex, in the hippocampus, to enable the formation of
453 detailed memories based on the specific combination of such invariant representations. This role
454 of the hippocampus aligns with our observation of mixed selectivity, i.e., neurons learning similar
455 feature combinations (Figure 5c), which in turn has been proposed to enable high-dimensional
456 representations important for higher cognition areas (Fusi et al., 2016; Bernardi et al., 2020).

457 **Circuit mechanisms underlying memory formation** The surprising observation of zero-shot
458 learning of place cells (Figure 4) suggests that the sparse autoencoders learned to cluster similar
459 samples onto single neurons in a generic manner. We hypothesize that the responsible circuits might
460 correspond to known hippocampal processes, mainly pattern completion and pattern separation (Rolls,
461 2013). On the one hand, neurons are pushed to collapse across-sample variability by clustering
462 samples based on similarity, an effect akin to pattern completion. On the other hand, sparsity also
463 imposes sharp discontinuities between clusters in neuronal space, even when they might be close in
464 input space, a process akin to pattern separation. The combination of both processes is reminiscent
465 of the locality-sensitive hashing (LSH) algorithms used in the computer science field for fast image
466 search (Kulis & Grauman, 2009). Future work will shed light on the learned circuit mechanisms
467 behind such a LSH in sparse autoencoders, and their potential mapping to pattern separation and
468 completion.

469 **Place cell over-representation near reward areas** We have shown that the development of
470 localized receptive fields can be generalized to other input modalities, such as sound. Crucially, we
471 used this simplified framework to demonstrate that receptive field distribution tends to be modulated
472 by the input sampling distribution (Figure 5d). Importantly, it has been observed that the density of
473 place fields increases near reward areas (Mamad et al., 2017). This has led researchers to seek external
474 reward signals in the hippocampus that would modulate the place cell map (Kaufman et al., 2020).
475 However, based on our results, we propose an alternative explanation based on oversampling: animals
476 tend to spend more time within rewarded areas due to consummatory behaviors, hence biasing sensory
477 sampling and learning. Additionally, in line with this hypothesis, place cell trajectories leading to
478 rewards or goals tend to be replayed more often than unrewarded past trajectories (Ambrose et al.,
479 2016), which would further reinforce the sampling bias.

480 **Biological plausibility** Although it has been claimed that error backpropagation and gradient
481 descent are mechanisms that could be implemented in the brain (Lillicrap et al., 2020), particularly in
482 the hippocampus (Santos-Pata et al., 2021b), we believe that such strong assumptions are unnecessary
483 to map our model and observations to the real hippocampus. The orthonormal activity regularization
484 term used in our sparse autoencoders could be realized by combining strong lateral inhibition (pro-
485 moting pairwise decorrelation) and homeostatic plasticity (ensuring that neurons maintain equalized
firing rates over time). Additionally, the orthonormal term can be thought of as sparse whitening

in ReLU-like neurons (i.e., pairwise decorrelation and variance normalization in the low-firing rate regime), a mechanism proposed to be realized in the brain by an overcomplete basis of inhibitory interneuron projections (Duong et al., 2023b;a; Lipshutz et al., 2022). Therefore, we contemplate several alternative mechanisms whereby the main objective driving our sparse autoencoders could be realized in brain circuits.

Limitations While our reinforcement learning experiments suggest that DQNs can make use of very high-dimensional representations to solve complex tasks, the present study is limited in scope (Figure 6). We tested only a few tasks (four tasks within the Animal-AI testbed) and a single model (DQN). To gain a more comprehensive understanding of the suitability of hippocampal-like representations for behavioral and policy learning, further testing is required with a broader range of tasks and models, especially in non-stationary environments where unseen samples are the norm. Additionally, future research should explore how these sparse autoencoders could enhance reinforcement learning algorithms that rely on discrete representations, potentially enabling algorithms based on, e.g., the successor representations (Dayan, 1993), to extend beyond simplified grid worlds.

REFERENCES

- R Ellen Ambrose, Brad E Pfeiffer, and David J Foster. Reverse replay of hippocampal place cells is uniquely modulated by changing reward. *Neuron*, 91(5):1124–1136, 2016.
- Dmitriy Aronov, Rhino Nevers, and David W Tank. Mapping of a non-spatial dimension by the hippocampal–entorhinal circuit. *Nature*, 543(7647):719–722, 2017.
- Omri Barak, Mattia Rigotti, and Stefano Fusi. The sparseness of mixed selectivity neurons controls the generalization–discrimination trade-off. *Journal of Neuroscience*, 33(9):3844–3856, 2013.
- C Bradford Barber, David P Dobkin, and Hannu Huhdanpaa. The quickhull algorithm for convex hulls. *ACM Transactions on Mathematical Software (TOMS)*, 22(4):469–483, 1996.
- Horace B Barlow et al. Possible principles underlying the transformation of sensory messages. *Sensory communication*, 1(01):217–233, 1961.
- Marcus K Benna and Stefano Fusi. Place cells may simply be memory cells: Memory compression leads to spatial tuning and history dependence. *Proceedings of the National Academy of Sciences*, 118(51):e2018422118, 2021.
- Silvia Bernardi, Marcus K Benna, Mattia Rigotti, Jérôme Munuera, Stefano Fusi, and C Daniel Salzman. The geometry of abstraction in the hippocampus and prefrontal cortex. *Cell*, 183(4):954–967, 2020.
- Benjamin Beyret, José Hernández-Orallo, Lucy Cheke, Marta Halina, Murray Shanahan, and Matthew Crosby. The animal-ai environment: Training and testing animal-like artificial cognition. *arXiv preprint arXiv:1909.07483*, 2019.
- Lara M Boyle, Lorenzo Posani, Sarah Irfan, Steven A Siegelbaum, and Stefano Fusi. Tuned geometries of hippocampal representations meet the computational demands of social memory. *Neuron*, 2024.
- Yusi Chen, Huanqiu Zhang, Mia Cameron, and Terrence Sejnowski. Predictive sequence learning in the hippocampal formation. *bioRxiv*, pp. 2022–05, 2022.
- Selmaan N Chettih, Emily L Mackevicius, Stephanie Hale, and Dmitriy Aronov. Barcoding of episodic memories in the hippocampus of a food-caching bird. *bioRxiv*, 2023.
- Peter Dayan. Improving generalization for temporal difference learning: The successor representation. *Neural computation*, 5(4):613–624, 1993.
- George Dragoi and György Buzsáki. Temporal encoding of place sequences by hippocampal cell assemblies. *Neuron*, 50(1):145–157, 2006.

- 540 Lyndon R Duong, David Lipshutz, David J Heeger, Dmitri B Chklovskii, and Eero P Simoncelli.
541 Statistical whitening of neural populations with gain-modulating interneurons. *arXiv preprint*
542 *arXiv:2301.11955*, 2023a.
- 543 Lyndon R Duong, Eero P Simoncelli, Dmitri B Chklovskii, and David Lipshutz. Adaptive whitening
544 with fast gain modulation and slow synaptic plasticity. *arXiv preprint arXiv:2308.13633*, 2023b.
- 545
546 Howard Eichenbaum. On the integration of space, time, and memory. *Neuron*, 95(5):1007–1018,
547 2017.
- 548 Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for
549 discovering clusters in large spatial databases with noise. In *KDD*, pp. 226–231, 1996.
- 550
551 André A Fenton and Robert U Muller. Place cell discharge is extremely variable during individual
552 passes of the rat through the firing field. *Proceedings of the National Academy of Sciences*, 95(6):
553 3182–3187, 1998.
- 554 Loren M Frank, Garrett B Stanley, and Emery N Brown. Hippocampal plasticity across multiple days
555 of exposure to novel environments. *Journal of Neuroscience*, 24(35):7681–7689, 2004.
- 556
557 Emmanouil Froudarakis, Uri Cohen, Maria Diamantaki, Edgar Y Walker, Jacob Reimer, Philipp
558 Berens, Haim Sompolinsky, and Andreas S Tolias. Object manifold geometry across the mouse
559 cortical visual hierarchy. *BioRxiv*, pp. 2020–08, 2020.
- 560 Stefano Fusi, Earl K Miller, and Mattia Rigotti. Why neurons mix: high dimensionality for higher
561 cognition. *Current opinion in neurobiology*, 37:66–74, 2016.
- 562
563 Sean Gillies. The shapely user manual. URL <https://pypi.org/project/Shapely>, 2013.
- 564
565 Mark A Gluck and Catherine E Myers. Hippocampal mediation of stimulus representation: A
566 computational theory. *Hippocampus*, 3(4):491–516, 1993.
- 567 James A Gornet and Matt Thomson. Automated construction of cognitive maps with predictive
568 coding. *bioRxiv*, pp. 2023–09, 2023.
- 569 Pablo E Jercog, Yashar Ahmadian, Caitlin Woodruff, Rajeev Deb-Sen, Laurence F Abbott, and Eric R
570 Kandel. Heading direction with respect to a reference point modulates place-cell activity. *Nature*
571 *communications*, 10(1):2333, 2019.
- 572
573 Alexandra Mansell Kaufman, Tristan Geiller, and Attila Losonczy. A role for the locus coeruleus
574 in hippocampal ca1 place cell reorganization during spatial reward learning. *Neuron*, 105(6):
575 1018–1026, 2020.
- 576 Nicholas Ketz, Srinimisha G Morkonda, and Randall C O’Reilly. Theta coordinated error-driven
577 learning in the hippocampus. *PLoS computational biology*, 9(6):e1003067, 2013.
- 578
579 Brian Kulis and Kristen Grauman. Kernelized locality-sensitive hashing for scalable image search.
580 In *2009 IEEE 12th international conference on computer vision*, pp. 2130–2137. IEEE, 2009.
- 581 Daniel Levenstein, Aleksei Efremov, Roy Henha Eyono, Adrien Peyrache, and Blake A Richards.
582 Sequential predictive learning is a unifying theory for hippocampal representation and replay.
583 *bioRxiv*, pp. 2024–04, 2024.
- 584
585 Michael S Lewicki and Terrence J Sejnowski. Learning overcomplete representations. *Neural*
586 *computation*, 12(2):337–365, 2000.
- 587
588 Timothy P Lillicrap, Adam Santoro, Luke Marris, Colin J Akerman, and Geoffrey Hinton. Backprop-
589 agation and the brain. *Nature Reviews Neuroscience*, 21(6):335–346, 2020.
- 590
591 David Lipshutz, Cengiz Pehlevan, and Dmitri B Chklovskii. Interneurons accelerate learning
592 dynamics in recurrent neural networks for statistical adaptation. *arXiv preprint arXiv:2209.10634*,
593 2022.
- 594
595 John Lisman. The theta/gamma discrete phase code occurring during the hippocampal phase precession
596 may be a more general brain coding scheme. *Hippocampus*, 15(7):913–922, 2005.

- 594 Lu Lu, Yeonjong Shin, Yanhui Su, and George Em Karniadakis. Dying relu and initialization: Theory
595 and numerical examples. *arXiv preprint arXiv:1903.06733*, 2019.
- 596
- 597 Omar Mamad, Lars Stumpp, Harold M McNamara, Charu Ramakrishnan, Karl Deisseroth, Richard B
598 Reilly, and Marian Tsanov. Place field assembly distribution encodes preferred locations. *PLoS*
599 *biology*, 15(9):e2002365, 2017.
- 600 Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and
601 projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- 602
- 603 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare,
604 Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control
605 through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 606 Edvard I Moser, Emilio Kropff, and May-Britt Moser. Place cells, grid cells, and the brain’s spatial
607 representation system. *Annu. Rev. Neurosci.*, 31:69–89, 2008.
- 608
- 609 May-Britt Moser, David C Rowland, and Edvard I Moser. Place cells, grid cells, and memory. *Cold*
610 *Spring Harbor perspectives in biology*, 7(2):a021808, 2015.
- 611 Josue Nassar, Piotr Sokol, SueYeon Chung, Kenneth D Harris, and Il Memming Park. On 1/n neural
612 representation and robustness. *Advances in Neural Information Processing Systems*, 33:6211–6222,
613 2020.
- 614 John O’Keefe and Jonathan Dostrovsky. The hippocampus as a spatial map: preliminary evidence
615 from unit activity in the freely-moving rat. *Brain research*, 1971.
- 616
- 617 Bruno A Olshausen and David J Field. Emergence of simple-cell receptive field properties by learning
618 a sparse code for natural images. *Nature*, 381(6583):607–609, 1996.
- 619 Shanshan Qin, Shiva Farashahi, David Lipshutz, Anirvan M Sengupta, Dmitri B Chklovskii, and
620 Cengiz Pehlevan. Coordinated drift of receptive fields in hebbian/anti-hebbian network models
621 during noisy representation learning. *Nature Neuroscience*, 26(2):339–349, 2023.
- 622
- 623 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
624 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
625 models from natural language supervision. In *International conference on machine learning*, pp.
626 8748–8763. PMLR, 2021.
- 627 Aviv Ratzon, Dori Derdikman, and Omri Barak. Representational drift as a result of implicit
628 regularization. *bioRxiv*, pp. 2023–05, 2023.
- 629
- 630 Stefano Recanatesi, Matthew Farrell, Guillaume Lajoie, Sophie Deneve, Mattia Rigotti, and Eric
631 Shea-Brown. Predictive learning as a network mechanism for extracting low-dimensional latent
632 space representations. *Nature communications*, 12(1):1417, 2021.
- 633
- 634 A David Redish. *Beyond the cognitive map: from place cells to episodic memory*. MIT press, 1999.
- 635
- 636 Edmund T Rolls. The mechanisms for pattern completion and pattern separation in the hippocampus.
637 *Frontiers in systems neuroscience*, 7:74, 2013.
- 638
- 639 Diogo Santos-Pata, Adrián F Amil, Ivan Georgiev Raikov, César Rennó-Costa, Anna Mura, Ivan
640 Soltész, and Paul FMJ Verschure. Entorhinal mismatch: A model of self-supervised learning in
641 the hippocampus. *Isience*, 24(4), 2021a.
- 642
- 643 Diogo Santos-Pata, Adrián F Amil, Ivan Georgiev Raikov, César Rennó-Costa, Anna Mura, Ivan
644 Soltész, and Paul FMJ Verschure. Epistemic autonomy: self-supervised learning in the mammalian
645 hippocampus. *Trends in Cognitive Sciences*, 25(7):582–595, 2021b.
- 646
- 647 Anirvan Sengupta, Cengiz Pehlevan, Mariano Tepper, Alexander Genkin, and Dmitri Chklovskii.
Manifold-tiling localized receptive fields are optimal in similarity-preserving neural networks.
Advances in neural information processing systems, 31, 2018.
- William Skaggs, Bruce McNaughton, and Katalin Gothard. An information-theoretic approach to
deciphering the hippocampal code. *Advances in neural information processing systems*, 5, 1992.

648 Carsen Stringer, Marius Pachitariu, Nicholas Steinmetz, Matteo Carandini, and Kenneth D Harris.
649 High-dimensional geometry of population responses in visual cortex. *Nature*, 571(7765):361–365,
650 2019.
651
652 Benigno Uria, Borja Ibarz, Andrea Banino, Vinicius Zambaldi, Dharshan Kumaran, Demis Hassabis,
653 Caswell Barry, and Charles Blundell. The spatial memory pipeline: a model of egocentric to
654 allocentric understanding in mammalian brains. *BioRxiv*, pp. 2020–11, 2020.
655
656 Reto Wyss, Peter König, and Paul FM J Verschure. A model of the ventral visual system based on
657 temporal stability and local memory. *PLoS biology*, 4(5):e120, 2006.
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A APPENDIX

A.1 DETAILED METHODS

A.1.1 MODEL’S ARCHITECTURE AND TRAINING

Visual autoencoder (Vis-AE) To compress the images from the Animal-AI environment, we employed a convolutional autoencoder (ConvAE) that maps the $3 \times 84 \times 84$ to a latent space with 1000 neurons, effectively compressing the images by a factor of ~ 21 . The details on the architecture can be found in Table 1.

Table 1: Architecture of the Visual Autoencoder (Vis-AE)

Layer	Type	Act. Func.	Filters/Units	Kernel Size	Stride/Padding
Input	Input	-	$3 \times 84 \times 84$	-	-
Conv1	2D Conv.	ReLU	16	4×4	2/1
Conv2	2D Conv.	ReLU	32	4×4	2/1
Conv3	2D Conv.	ReLU	64	4×4	2/1
Reshape	Reshape	-	6400	-	-
Fc1 (Z)	Linear	ReLU	1000	-	-
Fc2	Linear	ReLU	6400	-	-
Conv4	Trans. 2D Conv.	ReLU	32	4×4	2/1 (out. pad. 1)
Conv5	Trans. 2D Conv.	ReLU	16	4×4	2/1
Conv6	Trans. 2D Conv.	Sigmoid	3	4×4	2/1
Output	Output	-	$3 \times 84 \times 84$	-	-

Audio autoencoder (Aud-AE) To compress the frequency signals from the synthetic sound dataset, we employed another ConvAE that maps one-second time series data (with a resolution of 10^{-3} s) to a latent space with 100 neurons, effectively compressing the signals by a factor of 10. The details on the architecture can be found in Table 2.

Table 2: Architecture of the Audio Autoencoder (Aud-AE)

Layer	Type	Act. Func.	Filters/Units	Kernel Size	Stride/Padding
Input	Input	-	1×1000	-	-
Conv1	1D Conv.	ReLU	16	100	2/1
Conv2	1D Conv.	ReLU	32	100	2/1
Conv3	1D Conv.	ReLU	64	100	2/1
Reshape	Reshape	-	2624	-	-
Fc1 (Z)	Linear	ReLU	100	-	-
Fc2	Linear	ReLU	2624	-	-
Conv4	Trans. 1D Conv.	ReLU	32	100	2/1
Conv5	Trans. 1D Conv.	ReLU	16	100	2/1
Conv6	Trans. 1D Conv.	Tanh	1	100	2/1
Output	Output	-	1×1000	-	-

Loss function The loss function to be minimized in both autoencoders (Vis-AE and Aud-AE) includes a mean squared error (MSE) term as a reconstruction error to force the latent space to preserve the input information, and a orthonormal activity regularization term that promotes sparse representations in the latent space Z :

$$\mathcal{L} = \frac{1}{m} \|\mathbf{X} - \hat{\mathbf{X}}\|_2^2 + \frac{\lambda}{mn} \|\mathbf{I}_n - \mathbf{Z}^T \mathbf{Z}\|_F, \quad (2)$$

where \mathbf{X} is the input, $\hat{\mathbf{X}}$ is the output (i.e., the reconstructed input), λ is the regularization coefficient (10^3 by default), \mathbf{I}_n is the identity matrix with shape $n \times n$, \mathbf{Z} is the middle layer’s activity matrix of shape $m \times n$, m being the batch size, and n the number of hidden units. Therefore, the Gramian $\mathbf{Z}^T \mathbf{Z}$

(with shape $n \times n$) captures the pairwise co-activation strengths between neurons in latent space. The symbols $\|\cdot\|_2$ and $\|\cdot\|_F$ denote the squared L2-norm and the Frobenius norm, respectively. The orthonormal activity regularization term promotes pairwise neuron decorrelation while achieving an equalized contribution across neurons, alleviating the dying ReLU problem. For comparisons with the dense AE, λ was simply set to 0. We have found this orthonormal activity regularization to give improved and more reliable results than the L1 activity regularization term used in standard sparse autoencoders, especially in preventing the dead ReLU problem.

Data generation and training of Vis-AE For the visual experiments, datasets were generated by sampling a total of 10000 images in each of four Animal-AI environments: "doubleTmaze", "permanence", "cylinder", and "thorndike", at random locations within the arena (excluding the 10% of the space closest to each wall) and with random angles, following uniform distributions. Training was conducted using batches of 256 images for 10000 epochs, with independent training runs per environment. The Adam optimizer with no weight decay was used to train the network and the learning rate was set at 10^{-4} . In addition, the regularization strength λ was set to 10^3 for the sparse autoencoder, and 0 for the dense autoencoder. All weights were initialized using Xavier initialization except for Z (i.e., Fc1), whose weights were initialized following a random asymmetric initialization to minimize the dying ReLU problem (Lu et al., 2019). An early stop of 0.0005 in reconstruction loss was used to compare sparse and dense autoencoders with similar reconstruction capabilities.

Data generation and training of Aud-AE For the synthetic audio experiments, training consisted of batches of 256 one-second audio slices of varying frequencies. A sliding window of 1 second (with 1 ms shift) was applied to a linearly-varying frequency signal of total time 100 seconds, moving from 10 to 80 Hz, hence resulting in a total of 99001 samples. The sampling frequency was set at (10^4 Hz, so that the kernel size (1000) matched to one full cycle at the lowest input frequency (10 Hz). Training was conducted for 1000 epochs using the Adam optimizer with a learning rate of 10^{-4} , with no weight decay. Here, the regularization strength λ was set to 10^4 for the sparse autoencoder and 0 for the dense autoencoder. The weights were initialized as with the visual-AE, with Xavier initialization in all layers except for the the latent space Z that followed a random asymmetric initialization. An early stop of 0.002 in reconstruction loss was used to have a fair comparison between sparse and dense autoencoders.

A.1.2 SPATIAL TUNING

Firing ratemaps To generate ratemaps from latent space activity, we first created a grid of 60×60 bins (or 30×30 for computing spatial information scores) for each neuron. For each bin in the grid, we summed the neuron’s activity values for images sampled within that bin, generating an activity map in space. Then, an occupancy map was generated to account for the variability in the number of images sampled at each spatial bin (sampling density), which was used to normalize the values in the activity map. Finally, Gaussian smoothing was applied to each neuron’s normalized activity map, using a standard deviation of 3 bins. The resulting maps were normalized to their corresponding maximum values, yielding smooth ratemaps representing spatially-distributed neural activity.

Place field identification To identify and quantify place fields in each neuron’s ratemap, we first binarized the ratemap by setting pixels with activity below 20% of the maximum activity to zero (inactive bins) and those above to one (active bins). Clusters were identified by grouping adjacent active bins, forming a cluster if a group of active bins was completely surrounded by inactive bins. Clusters not meeting the size criteria for place cells (between 3% and 50% of the total number of bins, 3600) were discarded. The remaining clusters were considered place fields.

Spatial information Spatial information (SI) scores measure the amount of information a neuron’s firing rate (ν) conveys about the agent’s position (\mathbf{r}). For each neuron, we first normalized its ratemap (using a 30×30 bin grid) by the overall mean activity $\bar{\nu}$. Then, we computed an occupancy map that was normalized by the total number of samples to reflect the proportion of "time" spent in each bin of the ratemap, denoted as $p(\mathbf{r})$. Finally, we applied the formula introduced in Skaggs et al. (1992) to compute the SI scores:

$$SI = \sum_{\mathbf{r} \in \mathbf{R}} \frac{\nu(\mathbf{r})}{\bar{\nu}} \log_2 \left(\frac{\nu(\mathbf{r})}{\bar{\nu}} \right) p(\mathbf{r}). \quad (3)$$

The average SI across all neurons in the latent space Z provides an estimate of the degree of spatial tuning that the model has developed.

Spatial position decoding The spatial decoding error measures the expected error of a linear decoder using latent space activations Z to predict the spatial position \mathbf{r} . We fit a linear regression model with Z as the independent variables and \mathbf{r} as the dependent variables, predicting positions as $\hat{\mathbf{R}} = Z\mathbf{W}$. Then, we compute the mean squared error (MSE) between the predicted positions $\hat{\mathbf{R}}$ and the actual ones \mathbf{R} :

$$\text{MSE} = \frac{1}{n_{\text{samples}}} \|\mathbf{R} - \hat{\mathbf{R}}\|_2^2. \quad (4)$$

Finally, the average spatial decoding error (MSE) is re-scaled by dividing it by the maximum distance in the environment, that is, the diagonal of the arena, computed as $d = s\sqrt{2}$, with s being the side length.

A.1.3 INTERPRETABILITY

Visualizing and quantifying the network’s tiling of the image space We employed the CLIP neural network (Radford et al., 2021) to encode images (resized from $3 \times 84 \times 84$ to $3 \times 224 \times 224$) into 512-dimensional vectors. These vectors were subsequently reduced to a two-dimensional representation using UMAP (McInnes et al., 2018), with 10 neighbors and a minimum distance of 0.1, enabling the visualization of the high-dimensional image space.

Neurons in the hidden layer of the autoencoder that exhibited strong activation in response to specific images—those triggering activations exceeding a certain % of their maximum activation across the dataset—were mapped to points in the 2D image space. We then identify clusters of points using the DBSCAN algorithm (Ester et al., 1996), with radius ϵ of 1 and minimum samples of 4. These parameters are very dataset-dependent and were thus selected and validated via extensive visual inspection to ensure reliable cluster identification. Convex hulls were constructed around these clusters using the Quickhull algorithm (Barber et al., 1996) to delineate their spatial boundaries. This allowed us to identify the regions of the input space that each neuron encodes in their activations, i.e., their receptive fields.

Let $\{H_i\}$ denote the set of convex hulls corresponding to each neuron’s activated image space. The average overlap metric, \bar{O} , was calculated as follows:

$$\bar{O} = \frac{1}{\binom{k}{2}} \sum_{i < j} \frac{\text{Area}(H_i \cap H_j)}{\text{Area}(H_i \cup H_j)}, \quad (5)$$

where $\text{Area}(H_i \cap H_j)$ represents the area of intersection between hulls H_i and H_j , $\text{Area}(H_i \cup H_j)$ is the area of the union of hulls H_i and H_j , and k is the total number of hulls. The hull calculations were performed using the Shapely Python library (Gillies, 2013). The metric \bar{O} thus represents the average proportion of overlap relative to the union for each pair of hulls and ranges from 0 (no overlap) to 1 (complete overlap), thereby providing a quantitative measure of the redundancy in the neurons’ receptive fields across the image space.

Neuron clamping and decoding To test whether neurons in Z were directly interpretable based on their single-neuron activity (therefore obviating population codes), we conducted clamping experiments. This involved setting the activation of a specific neuron i in Z to its maximum activation value observed across the dataset \mathcal{X} , while setting the activations of all other neurons to zero. This is represented as $z'_i = (0, \dots, 0, x_{\max}, 0, \dots, 0)$ where $x_{\max} = \max(\{z_i | z = f(x), x \in \mathcal{X}\})$ and $f(x)$ represents the encoding function mapping \mathcal{X} to z . Then, z'_i is processed by the decoder $g(z'_i)$ (with $g(x)$ representing the decoding function mapping z to $\hat{\mathcal{X}}$) to yield an output signal (image of audio wave, depending on the AE).

Population code dimensionality The dimensionality of the population code was estimated by computing the power-law exponent α of the latent space activity Z (Stringer et al., 2019). We performed PCA on Z and computed the linear fit of the resulting eigenspectrum in log-log space over the range of the first 10 to 100 principal components. Since the exponent α provides an estimate of how fast the population activity eigenspectrum decays as new dimensions are added, high α values are indicative of low-dimensional codes, whereas low α values indicate high-dimensional codes.

864 A.1.4 REINFORCEMENT LEARNING EXPERIMENTS

865
866 **Animal-AI Testbed** The Animal-AI testbed is a comprehensive platform designed for evaluating
867 the cognitive and learning capabilities of AI agents in a variety of tasks that simulate real-world
868 challenges (Beyret et al., 2019). This testbed provides diverse environments where agents must use
869 visual cues and navigate complex structures to achieve specific goals. The visual inputs from these
870 environments are standardized to a resolution of 84 by 84 pixels, and agents can perform actions
871 defined by a 2-dimensional vector of integers: the first component goes from 0 to 2 and corresponds
872 to not moving, moving forward, or moving backwards, respectively; and the second component also
873 goes from 0 to 2 and corresponds to not rotate, rotate left, or rotate right, respectively. To encourage
874 efficient behavior, a standard frameskip of 4 is applied, and the reward value decreases by 0.001 at
875 each step. Episodes terminate either when the agent obtains the reward or after 1000 frames.

876 We evaluate our reinforcement learning agents using four distinct benchmarks within the Animal-AI
877 testbed: the Double T-maze, Object Permanence, Cylinder, and Thorndike tasks. Each of these tasks
878 presents unique challenges that require the agent to apply different strategies and cognitive abilities.

- 879 • **Double T-maze.** Each episode starts with the agent positioned randomly at one of the
880 corners of the maze, and the objective is to navigate to the center to obtain the reward. The
881 center contains the only positive reward (+3) available in the environment. Due to the high
882 and opaque maze walls, the agent cannot directly see the reward and must explore the maze
883 to find it.
- 884 • **Object Permanence.** At the beginning of each episode, the agent observes a large reward
885 (+3) falling behind a wall until it is completely occluded. The agent must then navigate to
886 the hidden reward, avoiding a small and visible reward (+1) along the way.
- 887 • **Cylinder.** This task involves an opaque cylinder with a medium-sized reward hidden inside.
888 The agent begins outside the cylinder and must navigate into the cylinder to obtain the
889 reward (+2).
- 890 • **Thorndike.** The task tests the agent’s ability to escape from a closed box to reach a reward
891 located outside the box. The box is semi-transparent, allowing the agent to see the reward
892 from inside. The only exit is blocked by a movable obstacle that the agent must push to
893 escape. A medium reward (+2) outside the box is the sole positive reward available.

894
895 **Model** To evaluate the performance of our sparse autoencoders in reinforcement learning scenarios,
896 we used a standard Deep Q-Network (DQN) architecture (Mnih et al., 2015) with modifications to
897 the input layer. Instead of feeding raw pixel data from the Animal-AI environments, we used the
898 compressed representations of 1000 units generated by the Visual Autoencoder (Vis-AE).

899 The loss function optimized by the Deep Q-Network (DQN) is the Mean Squared Error (MSE)
900 between the predicted Q-values and the target Q-values, calculated using the Bellman equation:

$$901 \mathcal{L}(\theta) = \mathbb{E}_{(s,a,r,s',d) \sim \text{ReplayBuffer}} \left[\left(r + \gamma \cdot (1 - d) \cdot \max_{a'} Q_{\text{target}}(s', a'; \theta^-) - Q_{\text{main}}(s, a; \theta) \right)^2 \right] \quad (6)$$

902 where Q_{main} is the main Q-network with parameters θ , Q_{target} is the target Q-network with parameters
903 θ^- , s is the current state, a is the action taken, r is the reward received, s' is the next state, d is a
904 boolean indicating whether s' is a terminal state, and γ is the discount factor. This loss function aims
905 to minimize the difference between the Q-value predicted by the main network and the target Q-value,
906 which is computed based on the reward and the maximum Q-value of the next state predicted by the
907 target network. The training of the DQN was performed by using the RMSprop optimizer. The target
908 network was periodically updated with the weights of the main DQN to stabilize training. The DQN
909 was trained with the following hyperparameters: a learning rate of 0.00025, a discount factor (γ) of
910 0.99, an update frequency of 4 steps, and a target network update frequency of 2500 steps. The ϵ for
911 the epsilon-greedy policy started at 1 and decayed linearly to 0.1 over 25000 steps. The replay buffer
912 size was set to 25000, with a batch size of 32 for experience replay. The details on the architecture
913 can be found in Table 3.

914
915 **Training and performance metrics** Each reported experiment tested two DQN agents, Sparse and
916 Dense, which differ only in their use of different Vis-AE models (sparse and dense autoencoders,
917 respectively) to obtain compressed representations from the environment observations as input. The

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

Table 3: Architecture of the Deep Q-Network (DQN)

Layer	Type	Act. Func.	Units
Input	Input	-	1000
Fc1	Linear	ReLU	100
Fc2	Linear	ReLU	50
Fc3	Linear	ReLU	25
Fc4	Linear	ReLU	9
Output	Output	-	9

two agents were evaluated across the four Animal-AI tasks described earlier. Each model run lasted 5000 episodes, and to ensure statistical reliability, each model played each task between 20 and 27 times. The reported average performance metric was calculated using a sliding window of 20 episodes.