# SELF-SUPERVISED GRID CELLS WITHOUT PATH INTE GRATION

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Grid cells, found in the medial Entorhinal Cortex, are known for their regular spatial firing patterns. These cells have been proposed as the neural solution to a range of computational tasks, from performing path integration, to serving as a metric for space. Their exact function, however, remains fiercely debated. In this work, we explore the consequences of demanding distance preservation over small spatial scales in networks subject to a capacity constraint. We consider two distinct self-supervised models, a feedforward network that learns to solve a purely spatial encoding task, and a recurrent network that solves the same problem during path integration. We find that this task leads to the emergence of highly grid cell-like representations in both networks. However, the recurrent network also features units with band-like representations. We subsequently prune velocity inputs to subsets of recurrent units, and find that their grid score is negatively correlated with path integration contribution. Thus, grid cells emerge without path integration in the feedforward network, and they appear substantially less important than band cells for path integration in the recurrent network. Our work provides a minimal model for learning grid-like spatial representations, and questions the role of grid cells as neural path integrators. Instead, it seems that distance preservation and high population capacity is a more likely candidate task for learning grid cells in artificial neural networks.

028 029

031

004

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

#### 1 INTRODUCTION

Known for their striking hexagonal spatial firing fields, grid cells (Hafting et al., 2005) of the medial Entorhinal Cortex (mEC) are thought to underpin several navigational abilities. These include path integration (Hafting et al., 2005; McNaughton et al., 2006; Burak and Fiete, 2009; Gil et al., 2017), forming a neural metric for space (Moser and Moser, 2008; Ginosar et al., 2023), vector navigation (Bush et al., 2015), and supporting memory and inference (Mulders et al., 2021; Whittington et al., 2020). Given the range of different functions believed to be supported by grid cells, it is natural to investigate which of these tasks, if any, are actually performed by these enigmatic cells.

In the modeling literature, emphasis has been placed on grid cells as path integrators, with compu-040 tational models establishing that grid cells are capable of doing path integration (Burak and Fiete, 041 2009). Recently, it has also been shown that grid-like representations *emerge* in neural networks 042 trained to path integrate (Cueva and Wei, 2018; Banino et al., 2018; Sorscher et al., 2022; Whitting-043 ton et al., 2020; Xu et al., 2022; Dorrell et al., 2022; Schaeffer et al., 2023), which has been taken 044 as evidence for grid cells performing path integration. However, under interventional cell ablations grid units are as important for path integration as randomly selected units (Nayebi et al., 2021) while band-like units are significantly more important (Schøyen et al., 2023). Moreover, these models are 046 typically complex, featuring different architectures and activation functions, interacting label cell 047 types (e.g. simulated place cell-like targets (Sorscher et al., 2022)), multiple regularization terms, 048 additional constraints (e.g. path invariance (Schaeffer et al., 2023) or conformal isometry (Xu et al., 2022)), and large cell counts. All of these works report grid-like representations, making it difficult to disentangle exactly what function grid cells serve in the various models. 051

In this work, we therefore propose a minimal model of grid cell function inspired by other recent models (Xu et al., 2022; Dorrell et al., 2022; Schaeffer et al., 2023), and use this model to approach the question of whether grid cells do path integration. Concretely, we are inspired by the notion

of grid cells serving as a metric for space, and consider an objective function that requires distance
 preservation over small spatial scales. In addition, we impose an L1 capacity constraint, which
 favors distributed representations that occupy a minimal portion of the state space. We train neural
 networks to minimize the proposed objective functions, and find that strikingly hexagonal grid-like
 spatial representations emerge using these two simple ingredients.

To explore whether grid cells do path integration in our model, we ablate path integration itself by training a feedforward (FF) network to minimize a purely spatial version of the proposed objective, alongside a recurrent neural network (RNN) tasked with implicit path integration. We find that the feedforward network learns grid representations on par with those of the path integrating RNN model. However, some RNN units display distinct band cell-like (Krupic et al., 2012) spatial responses, not seen in the FF network. When pruning velocity inputs to sampled subsets of units, we find that path integration contribution is inversely correlated with the mean sample grid score, with band-type cells providing the largest contribution.

Our findings suggest that grid cells may serve as a distributed high-capacity, distance-preserving representation. However, grid cells do not appear to be defined by the task of path integration. On the contrary, grid cells appear to be relatively unimportant for path integration, suggesting that grid cells may be more suitable for defining neural metrics for space, at least in artificial neural networks.

071 072

073 074

075

079

101 102 103

## 2 RESULTS & DISCUSSION

#### LOSS FUNCTION AND LEARNED REPRESENTATIONS

We consider the problem of training a representation that preserves distances in a neighborhood around a current location, as illustrated in Fig. 1a). Considering Cartesian coordinates  $\mathbf{x}_t$  (e.g. along a trajectory) where t indexes time, we propose the following objective

$$\mathcal{L} = \alpha \mathbb{E}_{t,t'} \left[ e^{-\frac{1}{2\sigma^2} \|\mathbf{x}_t - \mathbf{x}_{t'}\|^2} \left( \|\mathbf{x}_t - \mathbf{x}_{t'}\| - \|\mathbf{g}_t - \mathbf{g}_{t'}\| \right)^2 \right] + (1 - \alpha) \mathbb{E}[l_{cap}(\mathbf{g}_t)], \tag{1}$$

where  $\sigma$  is the envelope scale parameter that determines the width of the neighborhood distance preservation in the exponential term. This creates a window around each spatial location where the difference between physical and neural distances should be minimized. While this requirement is similar to demanding a conformal isometry (Xu et al., 2022), there are subtle, but behaviorally relevant differences (see Appendix A.9).

Gaussian radial basis functions are widely used and have been previously applied in e.g. normative 087 models of grid cells to promote local separation of neural representations, as seen in Dorrell et al. 880 (2022) and Schaeffer et al. (2023).  $\alpha$  is a hyper-parameter to weight the different loss terms, and 089  $\|\cdot\|$  denotes the Euclidean norm.  $\mathbf{g}_t$  is the representation we wish to learn, which we parametrize with either a feedforward neural network or a recurrent neural network, as illustrated in Fig. 1c) 091 and described in section 3.2. Although both models must solve the same spatial encoding task, the FF model takes direct Cartesian coordinate inputs, while the RNN only receives an initial Cartesian 092 position and subsequently receives Cartesian velocities. To correctly encode subsequent positions 093 and distances, the RNN therefore also needs to learn to path integrate. Notably, we constrain  $g_t$  to 094 be non-negative and of constant L2 norm, i.e.,  $g_{it} \ge 0$  and  $\|\mathbf{g}_t\| = 1$  for all i, t, similar to Xu et al. 095 (2022) and Schaeffer et al. (2023). The first loss term is minimized when Euclidean distances in the 096 learned representation g equal the target Euclidean distances in a neighborhood around the current 097 location. The second loss term,  $l_{cap}$ , is a capacity term. Xu et al. (2022) and Schaeffer et al. (2023) 098 posit that capacity constraints are conducive to grid-like representations, and Schaeffer et al. (2023) 099 proposed an L2 activity-based regularization term to maximize representational capacity. We, on 100 the other hand, propose using an L1 capacity term given by

$$l_{cap}(\mathbf{g}_t) = -\sum_i g_{it}.$$
(2)

Using an L1 capacity term is markedly different from an L2 capacity term both empirically and
mechanistically. An empirical investigation into the heterogeneous effects of using L2 capacity regularization instead of L1 is presented in Appendix A.8. A geometric illustration of when L1 capacity
is optimal is illustrated in Fig. 1b). When g has constant L2 norm and non-negative elements, L2
capacity promotes representations with similar angles, but any angle (in the positive quadrant) is

equally rewarded. Akin to L2 capacity constraints, the L1 capacity constraint (2) will also promote
 representations with similar directions. However, the L1 capacity constraint encourages maximally
 distributed and correlated cell activities. In other words, the full population vector state space is
 ideally placed near the diagonal vector, with all units coactive.

112 Surprisingly, both FF and RNN models learn highly grid-like representations, as seen in the ratemaps 113 in Fig. 1c), and quantified by grid scores in Fig. 1d). We further find in Fig. 1d) that their phase 114 distribution is seemingly random and uniform within the unit cell of the grid pattern. The histogram 115 of grid spacings is unimodally peaked, indicating a single module in both models. However, the 116 orientation histogram for the FF model is bimodal, suggesting two modules, but with identical spac-117 ing. In Appendix Fig. A1a) and b), we perform a parameter sweep across  $\alpha$  and  $\sigma$ , evaluating 118 grid score and grid spacing, and demonstrating how grid spacing can be tuned by adjusting these parameters. Moreover, Appendix Fig. A9 shows that grid spacing and field size vary independently 119 and can also be controlled using a third hyperparameter,  $\rho$ . Extending our model to include multiple 120 modules could be achieved by partitioning it, as in Xu et al. (2022), and assigning different distance-121 preservation hyperparameters to each module. However, we consider this beyond the scope of the 122 current study, as our focus is on analyzing the functional role of the emergent cell types while min-123 imizing potential confounding factors. Finally, Fig. 1e) shows the population vectors projected to 124 three dimensions using UMAP alongside persistence diagrams quantifying the number of persistent 125 1,2 and 3-dimensional cocycles. Both models show one 0D, two 1D, and one 2D hole, indicating 126 a toroidal manifold, which is also evident in the accompanying UMAP projection, consistent with 127 recent experimental evidence in biological grid cells (Gardner et al., 2022).

128 While both networks learn similar representations, we find that the RNN learned a small set of cells 129 with band-like representations, which is also visible in Fig. 1d) as a small bump near zero-grid score 130 (see Appendix A.6 for ratemaps of all units, for both networks). Why would the model add this extra 131 band-like subpopulation to the RNN? One key difference between the FF and the RNN is that the 132 RNN is required to path integrate. It therefore seems especially strange that the amount of grid-like 133 cells is reduced when the network is also required to path integrate, considering that mechanistic 134 theories advocate grid cells as the neural substrate for path integration, and normative models argue 135 that grid cells appear in RNNs trained for path integration. Additionally, when evaluating the loss terms on distinct subpopulations, we observe in Fig. A1b) that units with high grid scores achieve 136 a somewhat lower distance loss, while band-type cells afford slight improvements in capacity. This 137 supports the notion that grid cells are optimal for distance preservation. To rule out that architectural 138 differences induce the difference in observed patterns between model types, we also train RNNs 139 without velocity input, and find that band patterns vanish (see Appendix A.5). We also find that 140 band-type units become more prolific in networks trained in a high-speed setting (necessitating 141 stronger path integration).

- 142
- 143 144 145

#### 2.1 PRUNING & PATH INTEGRATION ABILITY

146 147

To investigate the role of different cell types during path integration, we selectively pruned velocity inputs to different cells as illustrated in Fig. 2a). By pruning velocity inputs, rather than network units directly (i.e. network states), we minimize off-target effects, as pattern formation and network stability should only be affected in cases where units require velocity input to perform state updates (i.e., path integrate). As an alternative to this approach, we also train a one-step linear decoder (see Appendix A7) to demonstrate that the network is in fact path integrating.

153 We categorized cell types with a grid score of less than 0.15 as band-like and the rest grid-like, 154 as seen in Fig. 2b). Pruning velocity input to the band-like cells induces a stark increase in (path 155 integration) error, as given by Eq. 4, over time, as shown in Fig. 2c). Comparably, pruning velocity 156 input to high grid score subpopulations showed almost no change in error over time. Furthermore, 157 we find that path integration performance is not strongly affected even when input to all grid units 158 in the network is pruned (see Fig. A2). Figure 2e) displays the corresponding error when pruning 159 uniformly across any cell (including both band and grid-like), which is accompanied by higher error. We hypothesise that this is due to the fact that low grid score units can be included in the 160 subpopulation. We also find that path integration ability correlates with band-like tuning, during 161 training (see Appendix A7).



193

194

197

199

Figure 1: Overview of models and objective function. a) Illustration of the objective function: distant locations should be represented by distant population vector, close locations by close population vectors. b) Illustration of the L1 capacity constraint. c) Illustration of the investigated neural net-195 work architectures (feedforward and recurrent) and inputs. Below, ratemaps of a random selection 196 of units are inset. d) Distributions of phases, grid scores, orientations, and spacing for both models. Orientations are given in radians, spacing is relative to environment dimensions (a  $4\pi \times 4\pi$  square arena). e) Persistence diagram and 3D UMAP projection of population ratemaps, for feedforward and recurrent networks. For the feedforward network, results are shown for units with orientation in the range [0.4, 0.8]. 200

- 201
- 202

203 Schøyen et al. found similar results when pruning band-like cells in Sorscher et al.'s RNN model of 204 grid cells. We further demonstrate in Fig. 2d) how the error is linearly related to grid score, where 205 pruning low grid score units has a high impact, and conversely, pruning high grid score units yields 206 low error. Finally, we also compare the initial state difference (ISD), as described in Equation 5, to 207 later states along different trajectories when pruning in Fig. 2f). We see the ISD rise fastest with no pruning, as would be expected when moving away from the initial state. When pruning high grid 208 score units, the trend is similar to no pruning, i.e., the neural representation is moved away from 209 its initial state over time, as one would expect if the network was still path integrating. However, 210 when pruning low grid score units (band cells), we see a much flatter ISD, indicating that the neural 211 state does not change as much. In other words, the neural path integrator is close to being turned 212 off. This provides compelling evidence that band cells, not grid cells, do path integration, at least in 213 recurrently connected artificial neural network models. 214

In other recent work that incorporates path integration in the learning task, most of the learned spatial 215 representations appear grid-like Xu et al. (2022); Schaeffer et al. (2023). Thus, in these models, it 216 appears that grid cells alone are responsible for the imposed objectives, including path integration. 217 Moreover, ablating each term of their losses provides compelling evidence for the need of each 218 component for robust pattern generation. However, Schaeffer et al. (2023) report band-like tunings 219 for some runs, and compared to our model, features a velocity-dependent, MLP recurrent weight 220 matrix that may obscure the contribution of other cell types (see Appendix A.4 for a detailed model comparison). In Xu et al. (2022), while a high proportion of cells were classified as grid cells, other 221 types, such as band-like cells, can be observed in some reported LSTM units. 222

223 Importantly, neither model has explored whether path integration is a necessary condition for grid 224 pattern emergence. While this would be an intriguing test, it is unclear how these models could 225 be adapted to non-path-integration domains, as certain features they identify as essential for pattern 226 formation—such as path invariance and trajectory permutations—rely on path integration. For a comprehensive comparison between our model and others, including Xu et al. (2024); Dorrell et al. 227 (2022); Sorscher et al. (2022); Xu et al. (2022); Schaeffer et al. (2023), see Appendix A.4. 228

- 229
- 230 231

#### 2.2 PATTERN FORMATION, CONNECTIVITY, AND GENERALIZABILITY

232 To investigate whether all hexagonal patterns are created equal, we examined the pattern formation, 233 connectivity structure, and generalization capabilities of the two models. In the bottom row of Fig. 3a), a sorted subset of feedforward unit ratemaps is shown, scaled by their outgoing weights to a 234 selected output unit (displayed as a large ratemap to the right), as detailed in section 3.5. For the FF 235 model, ratemaps of the penultimate layer serve as the basis representations forming the grid pattern 236 in the final layer. This basis representation is diverse, encompassing various patterns such as place-237 like (see the final three ratemaps) and single-band-like (see ratemap in the final row, second column) 238 ratemaps. These basis representations are non-periodic, indicating that the corresponding down-239 stream grid representations cannot maintain periodicity outside their training domain. This non-240 periodic nature can also be seen directly in the FF-network architecture, which uses non-periodic 241 activation functions. The bottom row of Fig. 3c) confirms that the network does not generalize the 242 grid pattern outside its training domain.

243 The RNN, while facing a similar out-of-domain initialization challenge as the FF network, can po-244 tentially generalize beyond the training arena boundaries using a learned, periodic path integrator 245 circuit. Subsequent ratemaps and their cumulative pattern formation rely on previously similar pe-246 riodic patterns, as observed in the top rows of Fig. 3a). Fig. 3c) demonstrates that initializing the 247 RNN inside its training domain and allowing it to path integrate along long sequences beyond the do-248 main reveals a clear periodic generalization. Conversely, starting outside its training domain results 249 in the same non-generalization behavior as the FF network. Additionally, Fig. 3b) shows that the 250 connectivity profile of the recurrent cells follows a short-range excitation and long-range inhibition 251 principle with respect to cells with neighboring phases, which is a known structure for generating grid-like representations (Burak and Fiete, 2009). Interestingly, the band-like cells exhibit an 252 excitation-inhibition connectivity profile shifted by their phase. Combined with the previous finding 253 that band cells function as the neural path integrator circuit in the RNN, this suggests a mechanism 254 for path integration through an excitation-inhibition shift of population activity in the wave direction 255 of the band. Fig. 3b) also shows that the input weight matrix has a particular structure: velocities 256 are projected along band directions (forming a hexagonal pattern in connectivity). Furthermore, 257 velocity weights tend to decrease with increasing grid score, possibly suggesting that grid cells are 258 less tuned towards integrating velocities.

259

261

#### 260 2.3 SUMMARY AND OUTLOOK

262 Our approach, featuring a minimal model (only two loss terms and a model without path integration), 263 has allowed us to isolate key factors contributing to the emergence of grid-like representations. 264 This simplification contrasts with previous complex models, which, while effective in generating 265 grid-like patterns, often obscure the specific contributions of grid cells to navigational tasks due to 266 their many-sided nature. Our findings demonstrate that grid cells, as encoded in feed-forward and 267 recurrent neural networks, can emerge when optimized to preserve local spatial distances under a capacity constraint. This aligns with previous theories positing that grid cells contribute to spatial 268 navigation by providing a consistent metric for space (Moser and Moser, 2008; Stemmler et al., 269 2015; Ginosar et al., 2023; Schøyen et al., 2024).



Figure 2: Results of pruning velocity inputs to the RNN. a) Illustration of the velocity input 309 pruning: during pruning, velocity projections to randomly selected subsets of units are silenced. 310 b) Grid score distribution, with indicated cutoff for low grid score units. Shown are also ratemaps 311 for low and high grid score units. c) Path integration error for pruning of 1000 subpopulations 312 among units with high grid score. Also shown is the error for low grid score units (dashed). d) Path 313 integration error at each trajectory timestep, for pruning of randomly sampled subpopulations of 314 units. Inset is the Pearson correlation coefficient between subpopulation mean grid score and the PI 315 error. e) As in c), for random subsamplings of the full population. f) Average initial state distance, for pruning random subpopulations of the full population. Also shown is the initial state distance 316 for the low grid score selection (dashed). 317

- 318
- 319
- 320
- 321
- 322
- 323



345 Figure 3: Pattern formation and out-of-bounds behaviors. a) Example pattern formation for 346 select recurrent units (top and middle), and a feedforward unit (bottom); ratemaps inset on the 347 right.  $g_{prev}$  denotes previous activations (penultimate layer activation for FF, second-to-last state for RNN), while cumulative representations shows the cumulative sum of weighted unit inputs. b) 348 Top: Ratemaps, alongside spatial connectivity of the recurrent network. Shown is the outgoing 349 recurrent weight, as a function of unit spatial phase. Also inset is the convex hull of all ensemble 350 phases. Red indicates excitation (positive weight), blue inhibition (negative weight). Bottom: x- and 351 y- components of the input matrix of the RNN, shaded by grid score of target unit. c) Evaluation 352 of recurrent (top) and feedforward units (bottom) units when the network is the environment is 353 extended beyond the original training regime (white square). For the RNN, representations are 354 shown for trajectories starting within the training enclosure, and both inside and outside.

- 355
- 356 357
- 358
- 359

360

Our results challenge the widely-held view that grid cells are integral to path integration. In our experiments, feed-forward networks, which lack an explicit path integration mechanism, still developed grid-like representations comparable to those in the path-integrating recurrent network. This observation suggests that the emergence of grid cells is not contingent on the process of path integration. Instead, it indicates that grid cells may function to encode spatial relationships.

Pruning velocity inputs further allowed us to disentangle the contributions of different cell types
to path integration. By comparing the representations of a pruned network to a non-pruned network, we could directly ascertain that path integration was not critically dependent on grid cells,
but seemingly rather on band-type units, echoing other recent findings (Schøyen et al., 2023). The
inverse correlation between the contribution to path integration and the mean sample grid score further reinforces the notion that grid cells might not be as crucial for path integration as previously
thought.

In summary, our findings suggests a reevaluation of the role of grid cells within the neural navigation
framework. While grid cells clearly provide a powerful spatial metric, their necessity for path integration is less obvious. This insight not only advances our understanding of grid cell functionality
but also prompts a reconsideration of how spatial representations are constructed and utilized in both
biological and artificial systems. Future research should continue to explore the interplay between
different cell types in the entorhinal cortex, as well as investigate how these findings can inform the
development of more efficient and interpretable models of spatial navigation.

# 378 2.4 LIMITATIONS 379

While our work provides evidence supporting the role of grid cells as a high capacity distance preserving representation, our model operates in a simplified domain, with Cartesian coordinate and velocity inputs. One could therefore consider extending the recurrent model to include e.g. more expressive input projections (such as Schaeffer et al. (2023)), or use other self-motion information, such as simulated head direction cell (Taube et al., 1990) input.

385 From a biological perspective, the use of label distances during training is also implausible. A related 386 limitation is the use of Euclidean distances in computing the loss function. While this might be 387 sufficient in an open arena such the one used in this work, more naturalistic settings, with e.g. interior walls may require other distance functions to capture the observed deformation of grid patterns in 388 e.g. exotic geometries (Ginosar et al., 2023). Using the Euclidean distance between state vectors 389 may also lead to inaccuracies when computing the loss. As an alternative, one could compute state 390 distances using the metric induced by the neural network model. However, when comparing over 391 smaller distances (which are implicitly enforced by  $\sigma$ ), the Euclidean distance function could still 392 be a fair approximation. 393

394

## 3 Methods

395 396 397

398

#### 3.1 DATA & INPUT TO NETWORKS

We considered two distinct neural networks in this work, with different inputs. The feedforward network received batches of Cartesian coordinate inputs (x, y), sampled randomly and uniformly from a square region with side lengths  $4\pi \times 4\pi$  (arbitrary units). The recurrent network, on the other hand, received Cartesian velocity inputs  $(v_x, v_y)$ , along trajectories sampled in the same square region. Additionally, the recurrent network was given the starting location  $(x_0, y_0)$  of each trajectory (in Cartesian coordinates), to form a suitable initial state.

To generate trajectories, a bouncing procedure was used. Starting locations were sampled randomly 405 and uniformly within the square arena. At each step, head directions were sampled according to 406 a von Mises distribution with scale  $\kappa = 4\pi$ , and step sizes drawn from a Rayleigh distribution 407 with scale parameter s = 0.15. Subsequently, we checked whether the resulting step landed out-408 side the enclosure. If so, the component of the velocity vector normal to the offending boundary 409 was reversed, effectively causing an elastic collision with the wall, keeping the trajectory inside. 410 Otherwise, the procedure was repeated until the desired number of timesteps and trajectories was 411 achieved. Due to the simplicity of the data, no datasets were pre-created, and all training data was 412 new to the networks, i.e. created on the fly.

413 414

415

## 3.2 NEURAL NETWORKS & TRAINING

416 We consider two distinct architectures: A fully connected feedforward network and a recurrent neural network. The FF model consisted of two hidden layers, with 64 and 128 units, respectively, 417 followed by an output layer of size  $n_q = 256$  units, transforming 2D Cartesian coordinates to a 418 latent space of  $n_q$  dimensions. We applied the ReLU activation function after each hidden layer, and 419 normReLU after the output layer. normReLU is a normalized ReLU function, which we take to be 420 given by normReLU( $\mathbf{x}$ ) = ReLU( $\mathbf{x}$ )/maximum(||ReLU( $\mathbf{x}$ )||,  $\varepsilon$ ), with  $\varepsilon = 10^{-12}$  a small constant 421 to avoid zero division. Finally, we initialize the weights of each layer in the FF model uniformly 422 between  $-\sqrt{k}$  and  $\sqrt{k}$ , where k is the number of input features to the layer. 423

Similar to the feedforward model, the RNN model featured  $n_g = 256$  recurrently connected units. The state of the RNN model at a time t was given by

426 427

$$\mathbf{g}_t = \operatorname{normReLU}(W\mathbf{g}_{t-1} + W_{in}\mathbf{v}_t)$$
(3)

where  $W \in \mathbb{R}^{n_g \times n_g}$  is a matrix of recurrent weights,  $W_{in} \in \mathbb{R}^{n_g \times 2}$  a matrix of input weights, while  $\mathbf{v}_t$  the velocity input at t. The initial state of the RNN was encoded with a feedforward neural network, with the exact same architecture as the FF model. In this case, the initial Cartesian position of the intended trajectory was provided as input. The weights in the recurrent matrix Wwere initialized to the identity to mitigate vanishing/exploding gradients, similar to Le et al. (2015), while the weights of the input matrix  $W_{in}$  were initialized uniformly, similarly to the FF model. We also initialized the RNN according to a uniform distribution (Fig. A3), with similar results, but slower convergence times. We therefore use the identity initialization throughout this work.

Both models were trained with a mini-batch size of 64 using the Adam optimizer with a learning rate of  $10^{-3}$  (Kingma and Ba, 2017). We set  $n_g = 256$  for both models, and trained the RNN on 10-timestep trajectories. The RNN was trained for a total of 50000 training steps and the FF network for 100000 steps, on unseen data. Note that the RNN initial state encoder was trained from scratch when training the RNN, and did not reuse the previous, independent FF model.

The final parameters needed for training are  $\sigma$  and  $\alpha$ , which define the loss function (1).  $\sigma$  is an 441 envelope scale parameter for the distance preservation loss term (See Appendix A.9 for more).  $\alpha$ 442 determines the relative weighting of the two terms in the loss, where  $\alpha$  is the weight for the distance 443 loss term and  $1 - \alpha$  is the weight for the capacity loss term, chosen such that  $\alpha \in [0, 1]$ . We 444 performed a grid search for  $\alpha$  and  $\sigma$  for both models to explore the impact of this weighting on the 445 grid score and scale of the representations. The values for  $\alpha$  were chosen uniformly between 0.01 446 and 0.99, while  $\sigma$  values were selected empirically around values that were seen to yield models 447 high grid score during initial testing. After the grid search (Fig. A1 a)), the values  $\sigma = 1.2$  and 448  $\alpha = 0.54$  were chosen for the FF and RNN models as this yielded high grid scores for both models, 449 even though multiple other combinations of  $\alpha$  and  $\sigma$  gave similar results. See Appendix A.1 and A.6 for more on the effects of  $\alpha$  on learned representations. 450

To investigate the influence of different RNN unit types on path integration ability, we ablated velocity inputs to subsets of recurrent units. Specifically, the velocity-pruned recurrent state was given by

456 457

468

469 470

471

472 473 474

475

451

453

$$\tilde{\mathbf{g}}_t = \operatorname{normReLU}(W\mathbf{g}_{t-1} + \mathbf{m} \odot (W_{in}\mathbf{v}_t)),$$

where **m** is a binary mask that silences velocity input to select units, and  $\odot$  the Hadamard product. RNN units were subselected so that a given subpopulation featured as many units as the number of low-gridscore units (n = 29; grid score cutoff 0.15). The subselection procedure was done for the full ensemble (random), as well as high grid-score units (grid score above cutoff). In each case, 1000 subpopulations were sampled randomly, with equal probability across units.

Since the proposed objective does not feature any decoding into a known target representation such as Euclidean coordinates, an alternate method of assessing whether the network is path integrating correctly is needed. We therefore computed the mean square error between velocity-pruned representations and a target representation created by running the network on the same velocity input without pruning. In other words,

$$\operatorname{Error}_{t} = \frac{1}{M} \sum_{i=1}^{M} ||\mathbf{g}_{t}^{i} - \tilde{\mathbf{g}}_{t}^{i}||^{2}, \tag{4}$$

where M = 10000 is the number of evaluated trajectories, and T = 10 the trajectory length. As a baseline, we also computed the mean squared difference with the initial RNN state, i.e.

$$ISD_{t} = \frac{1}{M} \sum_{i=1}^{M} ||\mathbf{g}_{0}^{i} - \tilde{\mathbf{g}}_{t}^{i}||^{2},$$
(5)

for velocity pruning to randomly sampled subpopulations, as described previously. To further establish that the network is capable of path integration, we also explore a trainable, one-step linear decoder (See. Appendix A.7).

# 480 3.4 GRID STATISTICS 481

To quantify the properties of the learned representations, we compute the grid score, orientation, spacing and phase of smoothed unit ratemaps. Smoothing was done with a Gaussian kernel, using the Astropy library (The Astropy Collaboration, 2022), with a standard deviation of 2 pixels, for  $64 \times 64$  ratemaps of unit activity. Grid scores were computed as the difference between the smallest and largest correlations, when comparing autocorrelogram annuli correlations at 60 and 30 degrees,

486 respectively. Orientations were computed as the smallest angle of the six innermost peaks of the 487 autocorrelogram (excluding the origin) with the horizontal. The grid spacing was computed as the 488 average distance to the six innermost peaks. Finally, grid phases were taken to be the displacement 489 of the closest grid peak to the origin of the ratemap.

490 491 492

493

494

495

496

#### 3.5 PATTERN FORMATION

To better understand how representations emerge in the feedforward and recurrent networks, we visualized the pattern formation process for select units. To do so, we weighted the activity of input ratemaps to a given unit by the network weight relating them. Then, weighted inputs were sorted according to their L2 norm over all spatial bins, and pattern forming was visualized using the cumulative sum of all weighted input ratemaps.

497 498 499

#### 3.6 OUT-OF-BOUNDS EVALUATION

500 501

To explore the ability of the networks to generalize beyond their training regime, we evaluated 502 trained feedforward and recurrent networks outside of the square training domain. Concretely, we 503 ran the feedforward network with Cartesian coordinates sampled in a  $8\pi \times 8\pi$  square, i.e., double 504 the original enclosure wall lengths. For the RNN, we evaluated two cases, one in which the network 505 started in the original square, and was allowed to move outside, and another wherein the network 506 started out anywhere in the enlarged enclosure. In both cases, the RNN was run for a total of 50 timesteps, and responses aggregated over 10000 trajectories.

508 510

511

507

#### 3.7 SUBPOPULATION LOSS EVALUATION

To assess whether different cell types contributed differently to the various loss terms, we evaluated 512 randomly sampled subpopulations of recurrent units on the capacity and distance losses separately. 513 For a given sample, the subpopulation vector consisted of the flattened ratemaps of the selected units. 514 Subsequently, the subpopulation vector was re-normalized to allow for a more direct comparison 515 with the full network loss. As a reference, we also evaluated low grid score units (n = 29; grid 516 score cutoff = 0.15) on either loss term. Additionally, we computed a baseline loss by randomly 517 shuffling the bins of unit ratemaps, effectively achieving a spatially random representation. All 518 subpopulations featured the same number of units (n = 29), as with the low grid score ensemble), 519 and a total of 1000 samplings were performed.

520 521

522 523

#### 3.8 PERSISTENT HOMOLOGY & LOW DIMENSIONAL PROJECTION

We used persistent homology to investigate whether the learned representation conformed to a sim-524 ple geometric structure. Specifically, we used the Ripser python library (Tralie et al., 2018) to com-525 pute persistence diagrams using spatially flattened ratemaps of unit activity. For the feedforward 526 network, the outermost 20 % of ratemap pixels were excluded to avoid boundary effects, and only 527 units with orientation in the range [0.4, 0.8] were included. This was done to avoid mixing possibly 528 independent modules of units with different orientations. For the RNN ratemaps, the ratemaps of 529 the full domain were used, and every unit was included in the analysis. For the analysis, we set 530  $n_perm = 500$ , and  $max_dim = 2$ , with otherwise default parameters.

531 We further visualized the neural structure of space using UMAP (McInnes et al., 2020) with 532  $n_neighbours = 2000$  and  $min_dist = 0.8$  to project ratemaps of the feedforward and recurrent networks down to three dimensions. We colored the resulting 3D point cloud using the first prin-534 ciple component of the ratemaps, similar to Gardner et al. (2022), highlighting important regions. 535 Ratemaps were processed in the same way as for the persistent homology analysis.

- 536 537
- 3.9 FIGURES

538 539

Figures were created using BioRender.com.

#### 540 REFERENCES 541

565

566

567

583

542	Andrea Banino, Caswell Barry, Benigno Uria, Charles Blundell, Timothy Lillicrap, Piotr Mirowski, Alexander
5/2	Pritzel, Martin J. Chadwick, Thomas Degris, Joseph Modayil, Greg Wayne, Hubert Soyer, Fabio Viola,
343	Brian Zhang, Ross Goroshin, Neil Rabinowitz, Razvan Pascanu, Charlie Beattie, Stig Petersen, Amir Sadik,
544	Stephen Gaffney, Helen King, Koray Kavukcuoglu, Demis Hassabis, Raia Hadsell, and Dharshan Kumaran.
545	Vector-Based Navigation Using Grid-like Representations in Artificial Agents. Nature, 557(7705):429-
546	433, May 2018. ISSN 0028-0836, 1476-4687. doi: 10.1038/s41586-018-0102-6. URL http://www.
547	nature.com/articles/s41586-018-0102-6.

- 548 Yoram Burak and Ila R. Fiete. Accurate Path Integration in Continuous Attractor Network Models of Grid Cells. 549 PLoS Computational Biology, 5(2):e1000291, February 2009. ISSN 1553-7358. doi: 10.1371/journal.pcbi. 550 1000291. URL https://dx.plos.org/10.1371/journal.pcbi.1000291.
- Daniel Bush, Caswell Barry, Daniel Manson, and Neil Burgess. Using Grid Cells for Navigation. Neuron, 552 87(3):507-520, August 2015. ISSN 08966273. doi: 10.1016/j.neuron.2015.07.006. URL https:// 553 linkinghub.elsevier.com/retrieve/pii/S0896627315006285. 554
- 555 Christopher J. Cueva and Xue-Xin Wei. Emergence of Grid-like Representations by Training Recurrent Neural Networks to Perform Spatial Localization. arXiv: 1803.07770, March 2018. URL http://arxiv.org/ 556 abs/1803.07770.
- 558 Yedidyah Dordek, Daniel Soudry, Ron Meir, and Dori Derdikman. Extracting Grid Cell Characteristics from 559 Place Cell Inputs Using Non-Negative Principal Component Analysis. eLife, 5:e10094, March 2016. ISSN 2050-084X. doi: 10.7554/eLife.10094. URL https://elifesciences.org/articles/10094. 560
- 561 William Dorrell, Peter E. Latham, Timothy E. J. Behrens, and James C. R. Whittington. Actionable Neural 562 Representations: Grid Cells from Minimal Constraints, September 2022. URL http://arxiv.org/ 563 abs/2209.15563. arXiv:2209.15563 [q-bio].
  - Ruiqi Gao, Jianwen Xie, Xue-Xin Wei, Song-Chun Zhu, and Ying Nian Wu. On Path Integration of Grid Cells: Group Representation and Isotropic Scaling. 2020. doi: 10.48550/ARXIV.2006.10259. URL https: //arxiv.org/abs/2006.10259.
- Richard J. Gardner, Erik Hermansen, Marius Pachitariu, Yoram Burak, Nils A. Baas, Benjamin A. Dunn, May-569 Britt Moser, and Edvard I. Moser. Toroidal Topology of Population Activity in Grid Cells. Nature, 602 (7895):123-128, February 2022. ISSN 0028-0836, 1476-4687. doi: 10.1038/s41586-021-04268-7. URL 570 https://www.nature.com/articles/s41586-021-04268-7. 571
- 572 Mariana Gil, Mihai Ancau, Magdalene I. Schlesiger, Angela Neitz, Kevin Allen, Rodrigo J. De Marco, and 573 Hannah Monyer. Impaired path integration in mice with disrupted grid cell firing. Nature Neuroscience, 574 21(1):81-91, December 2017. doi: 10.1038/s41593-017-0039-3. URL https://doi.org/10.1038/ s41593-017-0039-3. Publisher: Springer Science and Business Media LLC. 575
- Gily Ginosar, Johnatan Aljadeff, Liora Las, Dori Derdikman, and Nachum Ulanovsky. Are Grid Cells Used for 577 Navigation? On Local Metrics, Subjective Spaces, and Black Holes. Neuron, 2023. ISSN 0896-6273. doi: 578 10.1016/j.neuron.2023.03.027. URL https://www.sciencedirect.com/science/article/ pii/S0896627323002234.
- 580 Torkel Hafting, Marianne Fyhn, Sturla Molden, May-Britt Moser, and Edvard I. Moser. Microstructure of a 581 spatial map in the entorhinal cortex. Nature, 436(7052):801-806, 2005. ISSN 0028-0836, 1476-4687. doi: 582 10.1038/nature03721. URL http://www.nature.com/articles/nature03721.
- Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. arXiv: 1412.6980 [cs], 584 January 2017. URL http://arxiv.org/abs/1412.6980. 585
- Julija Krupic, Neil Burgess, and John O'Keefe. Neural Representations of Location Composed of Spatially Periodic Bands. Science, 337(6096):853-857, August 2012. ISSN 0036-8075, 1095-9203. doi: 10.1126/ 588 science.1222403. URL https://www.science.org/doi/10.1126/science.1222403.
- Quoc V. Le, Navdeep Jaitly, and Geoffrey E. Hinton. A Simple Way to Initialize Recurrent Networks of 590 Rectified Linear Units. arXiv: 1504.00941 [cs] Issue: arXiv:1504.00941 Publisher: arXiv, April 2015. URL 591 http://arxiv.org/abs/1504.00941. 592
- 593 Leland McInnes, John Healy, and James Melville. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction, 2020. \_eprint: 1802.03426.

- Bruce L. McNaughton, Francesco P. Battaglia, Ole Jensen, Edvard I Moser, and May-Britt Moser. Path integration and the neural basis of the 'cognitive map'. *Nature Reviews Neuroscience*, 7(8):663–678, August 2006.
  ISSN 1471-003X, 1471-0048. doi: 10.1038/nrn1932. URL https://www.nature.com/articles/ nrn1932.
- Edvard I. Moser and May-Britt Moser. A metric for space. *Hippocampus*, 18(12):1142–1156, December 2008.
   ISSN 1050-9631, 1098-1063. doi: 10.1002/hipo.20483. URL https://onlinelibrary.wiley.
   com/doi/10.1002/hipo.20483.
- Dounia Mulders, Man Yi Yim, Jae Sung Lee, Albert K. Lee, Thibaud Taillefumier, and Ila R. Fiete. A structured scaffold underlies activity in the hippocampus, November 2021. URL http://biorxiv.org/lookup/doi/10.1101/2021.11.20.469406.
- Aran Nayebi, Alexander Attinger, Malcolm Campbell, Kiah Hardcastle, Isabel Low, Caitlin S Mallory, Gabriel Mel, Ben Sorscher, Alex H Williams, Surya Ganguli, Lisa Giocomo, and Dan Yamins. Explaining heterogeneity in medial entorhinal cortex with task-driven neural networks. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. S. Liang, and J. Wortman Vaughan, editors, Advances in Neural Information Processing Systems, volume 34, pages 12167–12179. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper\_files/paper/2021/file/656f0dbf9392657eed7feefc486781fb-Paper.pdf.
- EunHye Park, Dino Dvorak, and André A. Fenton. Ensemble Place Codes in Hippocampus: CA1, CA3, and
  Dentate Gyrus Place Cells Have Multiple Place Fields in Large Environments. *PLoS ONE*, 6(7):e22349,
  July 2011. ISSN 1932-6203. doi: 10.1371/journal.pone.0022349. URL https://dx.plos.org/10.
  1371/journal.pone.0022349.
- William T. Redman, Francisco Acosta, Santiago Acosta-Mendoza, and Nina Miolane. Not so griddy: Internal representations of rnns path integrating more than one agent. *bioRxiv*, 2024. doi: 10.1101/2024.05.
  29.596500. URL https://www.biorxiv.org/content/early/2024/05/31/2024.05.29.
  596500.
- Rylan Schaeffer, Mikail Khona, Tzuhsuan Ma, Cristóbal Eyzaguirre, Sanmi Koyejo, and Ila Rani Fiete. Self-Supervised Learning of Representations for Space Generates Multi-Modular Grid Cells. 2023. doi: 10.
   48550/ARXIV.2311.02316. URL https://arxiv.org/abs/2311.02316. Publisher: arXiv Version Number: 1.
- Vemund Schøyen, Markus Borud Pettersen, Konstantin Holzhausen, Marianne Fyhn, Anders Malthe-Sørenssen, and Mikkel Elle Lepperød. Coherently remapping toroidal cells but not Grid cells are responsible for path integration in virtual agents. *iScience*, 26(11):108102, November 2023. ISSN 258900422. doi: 10.1016/j.isci.2023.108102. URL https://linkinghub.elsevier.com/retrieve/pii/ S258900422302179X.
- Vemund Schøyen, Constantin Bechkov, Markus Borud Pettersen, Erik Hermansen, Konstantin Holzhausen, Anders Malthe-Sørenssen, Marianne Fyhn, and Mikkel Elle Lepperød. Hexagons all the way down: Grid cells as a conformal isometric map of space. preprint, Neuroscience, February 2024. URL http:// biorxiv.org/lookup/doi/10.1101/2024.02.02.578585.
- Ben Sorscher, Gabriel C. Mel, Samuel A. Ocko, Lisa M. Giocomo, and Surya Ganguli. A unified theory for the computational and mechanistic origins of grid cells. *Neuron*, page S0896627322009072, October 2022. ISSN 08966273. doi: 10.1016/j.neuron.2022.10.003. URL https://linkinghub.elsevier.com/ retrieve/pii/S0896627322009072.

637

- Martin Stemmler, Alexander Mathis, and Andreas V. M. Herz. Connecting Multiple Spatial Scales to Decode the Population Activity of Grid Cells. *Science Advances*, 1(11):e1500816, December 2015. ISSN 2375-2548. doi: 10.1126/science.1500816. URL https://www.science.org/doi/10.1126/ science.1500816.
- Js Taube, Ru Muller, and Jb Ranck. Head-Direction Cells Recorded from the Postsubiculum in Freely Moving Rats. I. Description and Quantitative Analysis. *The Journal of Neuroscience*, 10(2):420–435, February 1990. ISSN 0270-6474, 1529-2401. doi: 10.1523/JNEUROSCI.10-02-00420.1990. URL https:// www.jneurosci.org/lookup/doi/10.1523/JNEUROSCI.10-02-00420.1990.
- The Astropy Collaboration. The Astropy Project: Sustaining and Growing a Community-oriented Open-source
   Project and the Latest Major Release (v5.0) of the Core Package, June 2022. URL http://arxiv.org/ abs/2206.14220. arXiv:2206.14220 [astro-ph].

648 649	Christopher Tralie, Nathaniel Saul, and Rann Bar-On. Ripser.py: A lean persistent homology library for python. <i>The Journal of Open Source Software</i> , 3(29):925, Sep 2018. doi: 10.21105/joss.00925. URL https:
650	//doi.org/10.21105/joss.00925.
651	James C.R. Whittington Timothy H. Muller, Shirley Mark, Guifen Chen, Caswell Barry, Neil Burgess
652	and Timothy E.J. Behrens. The Tolman-Eichenbaum Machine: Unifying Space and Relational Memory
653	through Generalization in the Hippocampal Formation. Cell, 183(5):1249-1263.e23, November 2020.
654	ISSN 00928674. doi: 10.1016/j.cell.2020.10.024. URL https://linkinghub.elsevier.com/
655	retrieve/pii/S009286742031388X.
656	Dehong Xu, Ruiqi Gao, Wen-Hao Zhang, Xue-Xin Wei, and Ying Nian Wu. Conformal Isometry of Lie Group
657	Representation in Recurrent Network of Grid Cells, 2022eprint: 2210.02684.
658	Dehong Xu, Ruigi Gao, Wen-Hao Zhang, Xue-Xin Wei, and Ying Nian Wu, Conformal Normalization in Re-
659	current Neural Network of Grid Cells, October 2023. URL http://arxiv.org/abs/2310.19192.
000	arXiv:2310.19192 [cs, q-bio, stat].
660	Debong Xu, Ruigi Gao, Wen-Hao, Zhang, Xue-Xin Wei, and Ying Nian Wu. An Investigation of Conformal
662	Isometry Hypothesis for Grid Cells, May 2024.
664	the set of
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	
600	
603	
694	
695	
696	
697	
698	
699	
700	
701	

#### A APPENDIX / SUPPLEMENTAL MATERIAL

#### A.1 HYPERPARAMETER SEARCH AND LOSS ABLATION

Fig. A1 shows the result of a hyperparameter search in the weighting coefficient  $\alpha$ , alongside training and evaluation loss for both models, as well as cell types. Note that  $\alpha = 0$  and  $\alpha = 1$  corresponds to complete ablation of distance and capacity losses, respectively.



Figure A1: Hyperparameter search, loss evaluation and training results. a) Grid score as a function of  $\sigma$  and  $\alpha$  for FF (top) and RNN models (bottom), where grid scale is calculated for the model with the highest grid score at a particular  $\sigma$ . b) Loss evaluation on 1000 randomly selected (random), spatially shuffled RNN unit ratemaps (random + shuffled) and low grid score (low GS) subpopulations. c) Loss training history for FF and RNN models.

749 750 751

702

703 704 705

706

707

708

709

## 752 A.2 EXTENDED PRUNING RESULTS

753

Fig. A2 shows ratemaps and error distributions for the recurrent network, in cases where velocity
 input to all grid and band units are ablated. Notably, path integration error is consistently lower
 when pruning input to grid than band units, even when every grid unit is velocity-deprived.



Figure A2: a) Ratemaps for every recurrent unit, after 10 timesteps, when pruning velocity inputs for: 200 grid units (left) and 29 band units (right). b) Final timestep error aggregated across 10000, 10-timestep trajectories as a function of number of pruned high grid score units. Shown is the median (black line) and the 25th and 75 percentile (shaded region. Also inset is a box plot of the corresponding error for pruning n = 29 low grid score units, with a corresponding jitter plot of the error distribution. Note the logarithmic scale.

Figure A3: Recurrent unit ratemaps for a network initialized according to (random) uniform distributions.

A.3 EFFECT OF NETWORK INITIALIZATION

Fig. A3 shows ratemaps of a recurrent network initialized according to a (random) uniform distribution, demonstrating that grid and band-like solutions are still learned without the identity initialization used elsewhere in this work.

A.4 COMPARISON TO OTHER MODELS

To contextualize our work, we compare it to several recent normative models of grid cells, highlighting both similarities and differences with the works Schaeffer et al. (2023); Sorscher et al. (2022); Dorrell et al. (2022); Dordek et al. (2016); Xu et al. (2022; 2024). Table 1 provides an overview of key characteristics of the models in Schaeffer et al. (2023); Sorscher et al. (2022); Xu et al. (2022), and our approach for both feedforward (FFN) and recurrent (RNN) implementations. It is especially noteworthy in the table that we are the only ones demonstrating emergence of grid-like cells with and without path integration (PI). Additionally, our model achieves high-fidelity grid patterns using a relatively simple architecture, requiring fewer loss components compared to Schaeffer et al. (2023) and omitting place-like encoding/decoding present in Sorscher et al. (2022) and Xu et al. (2022). 

Regarding emergence of grid-like cells without path integration, Dorrell et al. (2022) and Dordek et al. (2016) are worth highlighting. Dordek et al. (2016) demonstrates that grid cells can emerge as optimal low-dimensional, non-negative projections (non-negative PCA) of place-cell-like inputs modeled as difference-of-Gaussian radial basis functions. This framework elegantly connects place and grid cells, consistent with their known anatomical connections (Bush et al., 2015). However, this approach does not address normative aspects related to behaviorally relevant tasks, such as path integration or maintaining a distance-preserving representation. In contrast, Dorrell et al. (2022) employs a model constrained to a linear combination of sine and cosine plane waves. While grid-like patterns emerge, the simplicity of this formulation limits the flexibility of the resulting patterns compared to neural network-based approaches. Their functional constraints share similarities with the separation loss in Schaeffer et al. (2023), which encourages distinct neural representations for different spatial positions, but lack explicit distance-preserving objectives present in Xu et al. (2022), Schaeffer et al. (2023), and our model. 

860Xu et al. (2024) presents a simplified version of their earlier models (Xu et al., 2022; 2023; Gao et al.,8612020), removing the reliance on place-cell decoding that is prominent in Sorscher et al. (2022). This862refinement aligns more closely with Schaeffer et al. (2023) and our approach by imposing objectives863directly on grid representations  $g_t$ . However, their initial grid state is not explicitly stated complicat-863inject comparisons. Similar to Schaeffer et al. (2023) and our work, Xu et al. (2024) incorporates

864 a hard normalization constraint on grid representations. A key distinction between our FF model and 865 Xu et al. (2024) lies in the ablation of their second assumption—"transformation" (analogous to path 866 integration)-while still achieving grid-like cell emergence. This demonstrates the minimal require-867 ments necessary for producing grid-like cells, emphasizing the connection between grid cells and distance preservation rather than path integration. This is fundamental to our work, as we further 868 show that extending our FF model to an RNN framework, which incorporates path integration, leads 869 to the emergence of band-like representations. We subsequently provide extensive analysis showing 870 that these band-like representations are the primary functional component for path integration. 871

872 Whether Xu et al. (2024) also produce band-like cells is not clear. The authors report large en-873 semble grid scores, suggestive of all-grid representations. However, their models makes use of a 874 heading direction-dependent velocity input projection, which could conceivably produce band-like projections directly in the input. This architectural design is similar to Schaeffer et al. (2023), which 875 also reports high proportions of grid-like units (and some band cells, for certain runs), wherein the 876 recurrent weight matrix is an MLP and a function of the incoming velocity. Given the apparent 877 importance of band cells for path integration in our model, investigating whether these architectures 878 reproduce band-like behaviors could provide valuable insights into how differing model components 879 influence the emergence or suppression of band-like representations, offering a promising direction 880 for future research. 881

882

883

Table 1: Comparison of Grid Cell Models

4 Aspect 5	Schaeffer et al. (2023)	Xu et al. (2022)	Sorscher et al. (2022)	Ours FFN (+ RNN)
6 Objective	Sep + Inv + CI +	CI + PI	PI	Distance preser-
7	PI			vation (+ PI)
8 RNN Inputs	Learned RNN	Cartesian veloci-	Cartesian veloci-	None (+ Carte-
9	weight ma-	ties $\vec{v}_t$	ties $\vec{v}_t$	sian velocities
0	$\operatorname{trix}_{W(v_t)} W(v_t) =$			$v_t$ )
1	$MLP(v_t)$ from			
2	ties			
Initial state	N/A - the au-	Initial position in	Learned linear	Learned initial
Initial State	thors report that	place cell basis	projection of	representations
	the network is	$W_n \vec{p}(\vec{x}_0)$	initial position in	$MLP(\vec{x}_0)$
	initialized using		place cell basis	
	a "shared initial		$W_p \vec{p}(\vec{x}_0)$	
	state"			
PI architectu	e RNN-like with	Vanilla RNN	Vanilla RNN	None (+
	recurrent weight	and LSTM with	with ReLU	Normalized
	W(v) as input:	ReLU		vanilla RNN:
	W(v,),			Norm(ReLU( Wrd. + Wrd.))
Decoder	$(v_t)g_t)$	Linear readout to	Linear readout to	$W_Rg_t + W_Iv_t)$
Decouer	rione	place cell basis	place cell basis	ivone
		$\hat{\vec{n}} = W_{ij}\vec{a}_i$	$\hat{\vec{n}} = W_{\mu\nu}\vec{a}_{\nu}$	
Regularizatio	n L2 Capacity on	L2 weight	L2 weight	L1 capacity on $\vec{a}$
8	$\vec{q}_{t}$	penalty on Lin-	penalty on recur-	
	5.	ear place cell	rent weights $W_R$	
		readout weights		
		$W_{out}$		
Results	Emerges het-	Emerges high-	Emerges a range	Emerges only
	erogeneous	quality grid-like	of spatially	grid-like cells
	multi-modular	cells with	heterogeneous	with tunable
	grid-like cells	unable multi-	cells, including	multimodularity
	nlace and band	modulatily	griu-like cells	$(\tau )$ value-like
	like cells			
7				

#### 918 A.5 SPEED DEPENDENCE OF LEARNED REPRESENTATIONS 919

As we find that feedforward networks without velocity inputs learn grid-like representations, and that RNNs with velocity input learn additional band-like representations, it is natural to investigate the importance of the velocity signal for pattern formation in the RNN. We therefore trained recurrent networks on 10-timestep trajectory datasets with varying Rayleigh speed scales *s*. The results are presented in Fig. A4, which shows trained RNN unit ratemaps for s = 0, 0.075 and 0.75 (default value s = 0.15).

The case s = 0 corresponds to no speed, wherein the RNN sits idly at the starting location of each trajectory, not path integrating. Notably, patterns are still grid-like, but there is no apparent bandlike tuning. Thus, representations are similar to those of the feedforward network, and again, band structures appear linked to path integration. This also seems to indicate that band responses do not emerge as a consequence of architectural differences between networks, as part of a pattern forming process.

For small speeds (s = 0.075), some band-like tuning is again observed. However, for speeds much larger than that used in the default case (s = 0.75), a large fraction of units now appear strongly band-like. The pattern, surprisingly, has also turned square, rather than hexagonal, which has been observed in other path integrating models Cueva and Wei (2018). A square pattern may reflect that path integration is more accurately achieved along the cardinal directions of the input velocity vectors, but this observation warrants further investigation.

Together, these results all point towards bands being a necessary component of path integration
 in recurrent networks, as band-like representations emerge in proportion to the importance of path
 integration for the task at hand.

941 942

#### A.6 EXTENDED UNIT RATEMAPS

As noted previously, we find that feedforward networks learn grid-like representations, and that path integrating RNNs learn additional band-like structures. This is shown explicitly in Fig. A5, which shows every unit in two trained models, one feedforward and one recurrent. Notably, all appear grid-or band-like.

When ablating the capacity loss, however, RNN units tend to lose their regular firing patterns. In
particular, Fig. A6 shows that low-capacity units have heterogenous firing fields, reminiscent of
place fields, in large environments (Park et al. (2011)). While a minority of units do exhibit striped
or banded firing fields, none display clear periodic band-like tuning. One might therefore speculate
that path integrating band units are inextricably tied to grid cells, and that other cell types, such as
place cells, rely on distinct path integration mechanism.

954 955

964

#### A.7 DECODING & PATH INTEGRATION PERFORMANCE

Assessing path integration performance in the trained RNN is not straight-forward, due to the highly
periodic nature of its learned representations (which causes representations at different locations at
distant locations to be ambiguously encoded). However, a network may still be able to path integrate
correctly, even without global decodability (as we observe in our model, using the path integration
measures in Fig. 2). We therefore develop a simple, one-step linear decoder to uncover this ability,
and perform decoding *locally*.

- 962 Specifically, we consider a conditional linear decoding of the form
  - $\hat{\mathbf{x}}_t = W \mathbf{u}_t,$

where  $\mathbf{u}_t = \operatorname{cat}(\mathbf{g}_t, \hat{\mathbf{x}}_{t-1})$ , i.e. a concatenation of the current state  $\mathbf{g}_t$  and a previous estimated position  $\hat{\mathbf{x}}_{t-1}$ , and W is a trainable weight matrix.

967 We train this decoder in a one-step fashion by feeding in network states and (true) previous-step 968 locations  $x_t$ , along trajectories. However during inference, we decode positions iteratively by using 969 the previous position estimate  $\hat{x}_{t-1}$  produced by the decoder, so that the entire decoded trajectory is 970 produced by the decoder. To verify that the decoder is not simply copying the previous location, we 971 compare decodings of 1000-timestep, long-sequence trajectories starting at the origin, to a baseline 972 of always predicting the origin (which would be expected if the decoder simply copied its input).



1026					F	F					
1027	22 50							10 R			
1028	88 X.					0.5	S	<b>.</b>			
1029	**						8 8	8			86
1030	200 00			007	00 98		0 0	90.9			10
1031					<b>6</b> 5 26	200 4	2	202.0	<b>9</b> 900		
1032	** **	<b>X</b> 8					÷ 33	88 2			
1033	<b>44</b> 08	60 Q	x 304	20- 1		69 S	26 94	0.1	6 26	10	
1034		22.2		-02		02.2	0 00				
1035	<b>38</b> 38								0.91	<u> </u>	88
1036	84 88	<b>6</b> 8 0			<b>1</b>	02.8	8 00	Q. X		<b>6</b> 00	
1037	87 02			<b>***</b>						~	
1038	28 B)	S 8	88			38					
1039	<b>96 98</b>				<b>60</b> 602	<b>3</b> 2 3	8 50	90X 15	990		20
1040	20 20										
1041								88	5		
1042	22 22	X0 8		•••				30 8			
1044											
1045		88 A									
1046	80 SS	. ee S	•					<b>£</b>	8.89		38
1047						50 1					125
1048	<b>.</b>										
1049	88 88				<b>*</b> *		2 81	88 9			***
1050	50 89		8 30	20		0.0	5K 585				205
1051	<b>10</b> 00								е ж.		
1052	-				RI RI	NN			-		
1053						<b>6</b>		14 3	: ///		
1054						66 8	8		6		
1055											
1056			•	86						•••	
1057				83						83	60
1058	20 00										10
1059	<b>10</b>			88 E				-			77
1061	三次		3///				2 2	/// 3			
1062				69				00.5	8	111	
1063				22.5		222	0.10			111	
1064	$\parallel \parallel \equiv$	M 3			211			18 3			
1065	<b>60 92</b>					•••		<u> 22</u> 5			11
1066											
1067									8 19	38	
1068	9 B		3 ///				8.52				
1069											
1070								=;	10	111	<b>9</b>
1071	2 3	11					3 D		2 22	22	
1072			6 774								
1073			: 77		••	<b>.</b>	с н				
1074							// 2	2	1 S.		22
10/5	00 00	599 P			100		8	86 2	0 22		22
1077									•	•••	
1077											

Figure A5: Ratemaps of all 256 network units. Ratemaps of every network unit, for feedforward (top), and recurrent (bottom) networks.

Under review as a conference paper at ICLR 2025

1085																
1086																
1087																
1088																
1089																
1090																
1091																
1092																
1093			•	•	÷.	÷.		. ·		1	•	12	100	•	Ζ.	
1094					•	-	• •		•	•		20	٠.	•	<u> </u>	•
1095		•	٠.	•		14		٠.	٠.	14	150	•	1	<b>.</b>		•
1096				- 1 e	s. 1			•	•					1		
1097		<u>ر</u> ا		۰.	• 🔨	•		۰.	•	<u> </u>		•		1.	•	• *
1098	, * •		• • •	+ *		- * *	•	<u>.</u>	. · ·		•	1			•	1
1099				-	-	- · ·			-		•	-, '			•	
1100	- S	्_*	•		1	•			1		•	1.			, • ·	
1101	1	1		<ul> <li>1</li> </ul>	٠.	•		7.2	• •	10	- 11		۰.	•		
1102				<u> </u>	-		•		•		-		1	_	_	
1103		•					1		<b>.</b> .			×.,	×.	•••		
1104			~					•			•	•	•	•		
1105			Č 1	•	•				• 51		$\mathcal{I} = \mathcal{I}$	1	• *	1	٦.	•
1106		٠ <u>٢</u>		÷.,	•	$\mathcal{M}_{\mathcal{M}}$	•	• *	1			1	÷.,	٦٩.,	۰.	
1107			<i>*</i>			-			_	-	_					
1108	1.4		. 1			- 6-9	÷.,	•	1.	. •	1.	•	÷.	. •	•	11
1109			۰.	1	· •		~	• •		. <sup>1</sup>			1		5 m 1	1.1
1110		• • •	÷.,	× 1	• •		<u>`</u>		• •		•		- <b>s</b> *			÷.
1111	•	• -	11	۰.	÷.,	•	• 2	. •	•		. C	• •	÷.		. •	
1112		1	•		•	•	1				•••	ć.	•			
1113	· .	. A		÷.,	• •	• 1	۰.	•			$\mathbf{r}$	 	•	•	<u>, -</u>	17.7
1114	÷.,	• •	•	÷.,			٠,	• •	Č.	1	1		•	11	• •	•
1115						٠.					-		•			
1116		•	-	1		•		10		•						•
1117		•					•			٠.	•	-			۰.	٠
1118		• •		· •	٠.		•	1			8.0	٠			• •	
1119																

Figure A6: Low capacity networks learn irregular firing patterns. All recurrent units, for a network trained with  $\alpha = 1$  (no L1 capacity constraint).

The results are shown in Fig. A7a), which demonstrate that not only does the decoder outperform the zero-prediction baseline, but the network is capable of fairly accurate path integration for hundreds of timesteps. Notably, this extends far beyond the training sequence length of 10 timesteps, which is in line with the out-of-box generalization ability observed in Fig. 3c). The training loss in Fig. A7b) also demonstrate that the decoder becomes adept only after thousands of training steps. Fig. A7c) shows an example decoded long-sequence trajectory. Evidently, the decoded trajectory trails the true one for a long time, but shows signs of drift over time.

To explore how path integration ability changes over training time, we trained decoders at varying steps during model training (Fig. A7d)). As expected, path integration ability increases with training time. However, even an untrained network is somewhat decodable (possibly due to the inherent reservoir capacity of larger networks such as ours), but performance degrades considerably faster than trained networks. Another notable feature, is that fully trained networks do *not* appear to be the most capable path integrators, with the 1000-step checkpoint achieving the smallest decoding error.

Inspired by this, and our findings that band units are most important for path integration in trained networks, we examined network representations in terms of their *band score*. Following Redman et al. (2024), the band score is given by

1152

 $b(G) = \max_{k_x, k_y} \operatorname{Corrcoeff}(G, S(k_x, k_y)),$ 

1153 where  $k_x, k_y \in \{0, 0.1, 0.2, ..., 2\pi\}$  are spatial frequencies in the x and y directions, corrcoeff the 1154 Pearson correlation coefficient, and for G we use the autocorrelogram of the ratemap of a particular 1155 network unit (to ensure G is centered), and  $S(k_x, k_y)$  a 2D sinusoid whose frequency is given by  $k_x$ 1156 and  $k_y$ .

Intriguingly, we find that band scores indeed correlate with path integration ability (compare Fig. A7d) and e)). In particular, ensembles with greater band scores tend to perform better; compare e.g. the fully trained model, and the network trained for 1000 steps. As can be seen from inset ratemaps, ratemaps at this particular step are more square, and highly band-like.

As a final remark, this correspondence should not be taken to mean that path integration is necessarily performed by band units; it could for instance reflect that such units are more linearly decodable, for instance. However, these findings resonate with our other results (which do not explicitly rely on a linear decoder), again hinting at a connection between path integration ability and band-like representations.

- 1166
- 1167 1168

#### A.8 AN L2 CAPACITY CONSTRAINT INDUCES HETEROGENEOUS REPRESENTATIONS

In this work, we have shown that distance preservation and an L1 capacity constraint is sufficient to induce grid patterns in feedforward and recurrent networks. However, this choice was in part motivated by other work (Schaeffer et al., 2023), that utilizes a similar L2 constraint. To demonstrate why an L1 constraint, rather than an L2 appears more conducive to grid-like representations, we trained multiple networks with varying L2 capacity constraints.

1174 The resulting ratemaps are shown for both FF and RNN networks in Fig. A8. Notably, recurrent 1175 unit responses do appear hexagonal and grid-like (as well as band-like) for appropriate  $\alpha$  values. 1176 However, representations are more irregular than what we observe for an L1 constraint, with some 1177 units being mainly silent, some exhibit spurious firing fields, and others incomplete grid firing fields. 1178 Comparing to results for L1 capacity, grid cells are known for their tendency to fire all over the 1179 recording environment (in symmetric geometries), and exhibit persistent activity across different 1180 environments, suggesting an L1 constraint may be more appropriate.

1181 Unlike for the L1 case, we could not find hyperparameter values  $\alpha$  that resulted in grid-like repre-1182 sentations in the feedforward network using an L2 capacity constraint. Rather, the learned repre-1183 sentations tended to be sparse. Thus, an L1 capacity appears to favor distributed representations, 1184 aligning well with grid cell properties and potentially contributing to more robust representations. 1185 However, it should be pointed out that the presence of path integration induced band- and grid-like 1186 representations, also for the L2 constraint; suggesting that it may enhance pattern formation, and that 1187 band and grid patterns are better suited for simultaneous path integration and distance preservation than heterogeneous ones.



Figure A7: Recurrent network path integration performance. a) Path integration error (Eu-clidean distance between true and decoded trajectories) for the one-step linear decoder. Shown is the median error, for a trained RNN model, and a baseline, all-zero prediction, for 1000 trajectories all starting at the origin. Shaded regions indicate the inter-quartile range. b) Training loss for the one-step decoder. c) An example decoded true trajectories, alongside the corresponding true trajec-tory. d) Decoding error over trajectory time, as a function of model training length (legend indicates training step; 49999 denotes a fully trained RNN). e) Ensemble band scores corresponding to the training steps in d). Also inset are randomly selected unit ratemaps at the indicated training step. 



Figure A8: Effects of L2 capacity constraints on learned representations. Ratemaps of randomly selected units when the L1 capacity constraint is exchanged for an L2 constraint. Shown are responses for RNN (top) and FF (bottom) network units for varying loss weightings.  $\alpha = 0$ corresponds a pure capacity loss,  $\alpha = 1$  to a pure distance preservation loss.

1270 A.9 GRID SPACING AND FIRING FIELD SIZE

<sup>1272</sup> The grid search in A1 unveiled that the envelope scale parameter  $\sigma$  dictates the spacing of the grid <sup>1273</sup> pattern. This is shown explicitly in Fig. A9, wherein feedforward grid fields are shown to become <sup>1274</sup> more distantly spaced with increasing  $\sigma$ .

A related, interesting quantity is the grid firing field size. How this quantity is determined in our model, is not entirely obvious. However, because of the close correspondance between neural distances and physical ones, we note that we can modulate the scale of the pattern, by introducing a factor  $\rho$  into the loss function,

1263 1264

1269

$$\mathcal{L} = \alpha \mathbb{E}_{t,t'} \left[ e^{-\frac{1}{2\sigma^2} \|\mathbf{x}_t - \mathbf{x}_{t'}\|^2} \left( \rho \|\mathbf{x}_t - \mathbf{x}_{t'}\| - \|\mathbf{g}_t - \mathbf{g}_{t'}\| \right)^2 \right] + (1 - \alpha) \mathbb{E}[l_{cap}(\mathbf{g}_t)],$$

so that the representation learns to represent distances scaled by a factor of  $\rho$ . This is similar to the conformal scaling factor used by (Xu et al., 2022).

Ratemaps of feedforward units trained to minimize this slightly modified loss are also shown in Fig. A9, which demonstrates how firing fields become larger with decreasing  $\rho$ . Following (Xu et al., 2022) and (Xu et al., 2024), this control over learned representations allows us to directly introduce multiple modules in an interpretable way, by partitioning the network into distinct modules, or allowing  $\rho$  to be a unit-specific trainable parameter.

Thus, grid field size and spacing can be readily understood in our model: field sizes reflect the ratio between neural and physical distances, while grid spacing reflects the scale at which distances should be accurately represented (for a fixed field size).

This result also highlights an important, but subtle difference between distance preservation, and a conformal isometry requirement: While both enable faithful distance computations by integrating the (flat) metric along neural trajectories, demanding distance preservation allows distances to be computed directly by comparing two population vectors (while in the range where this is valid),





0.0 

Figure A10: Euclidean distances in representation vs. space. Distances between population vectors of feedforward networks trained with varying scale parameters,  $\sigma$ , versus distances in the arena. States and locations are sampled from a square grid, i.e. a ratemap. 

Physical Distance

 $\sigma = 3$ 

which could greatly simplify distance computations. However, it should be noted that this is also true to first order for a conformal isometry. 

That our model preserves Euclidean distances (and preservation is determined by  $\sigma$ ), is showcased in Fig. A10, where representational and physical Euclidean distances are compared directly. As shown, a larger value of  $\sigma$  induces a right-shift in the distance plot, and the relationship between the two is near-linear for longer, indicating that distances are preserved. 

Exploring the scale at which distances are preserved in biological grid cell data, could make for an interesting comparison between our model and others, and could possibly account for variations in grid spacing, which we find coincides with the scale at which (Euclidean) distances are preserved in the representation.