

No Free Lunch from Deep Learning in Neuroscience: A Case Study through Models of the Entorhinal-Hippocampal Circuit

Rylan Schaeffer^{1,2}, Mikail Khona², and Ila Rani Fiete^{2,3}

¹Computer Science, Stanford University

²Brain and Cognitive Sciences, Massachusetts Institute of Technology

³Physics, Massachusetts Institute of Technology

⁴McGovern Institute for Brain Research

August 7, 2022

Abstract

Research in Neuroscience, as in many scientific disciplines, is undergoing a renaissance based on deep learning. Unique to Neuroscience, deep learning models can be used not only as a tool but interpreted as models of the brain. The central claims of recent deep learning-based models of brain circuits are that they make novel predictions about neural phenomena or shed light on the fundamental functions being optimized. We show, through the case-study of grid cells in the entorhinal-hippocampal circuit, that one may get neither. We begin by reviewing the principles of grid cell mechanism and function obtained from first-principles modeling efforts, then rigorously examine the claims of deep learning models of grid cells. Using large-scale hyperparameter sweeps and theory-driven experimentation, we demonstrate that the results of such models may be more strongly driven by particular, non-fundamental, and post-hoc implementation choices than fundamental truths about neural circuits or the loss function(s) they might optimize. We discuss why these models cannot be expected to produce accurate models of the brain without the addition of substantial amounts of inductive bias, an informal No Free Lunch result for Neuroscience. Based on first principles work, we provide hypotheses for what additional loss functions will produce grid cells more robustly. In conclusion, caution and consideration, together with biological knowledge, are warranted in building and interpreting deep learning models in Neuroscience.

Introduction

Over the past decade, deep learning (DL) has underpinned nearly every success story in machine learning, e.g., [43, 4] and increasingly many advances in fundamental science research, e.g., [26]. In neuroscience, deep learning is similarly gaining widespread adoption as a useful method for behavioral and neural data analysis [40, 38, 21, 32, 30, 35].

But DL offers a unique contribution to neuroscience that goes beyond its role in other fields, in that deep networks can be viewed as models of the brain. The success of DL in matching or surpassing human performance means it is now possible to construct models of circuits that may underlie human intelligence. As a recent review wrote, “researchers are excited by the possibility that deep neural networks may offer theories of perception, cognition and action for biological brains. This approach has the potential to radically reshape our approach to understanding neural systems” [41]. Further, DL is a democratizing force for building neural circuit models of complex function.

The essential claims (and promises) of DL-based models of the brain are that 1) Because the models are trained on a specific optimization problem, if the resulting representations match what has been observed in the brain, then they reveal which optimization problem(s) the brain is solving, or 2) These models, when trained on sensible optimization problems, should generate novel predictions about the brain’s representations and behavior.

However, given the nascent nature of such approaches and the excitement accompanying some claims, we should examine them carefully. In DL and deep reinforcement learning, some successes attributed to

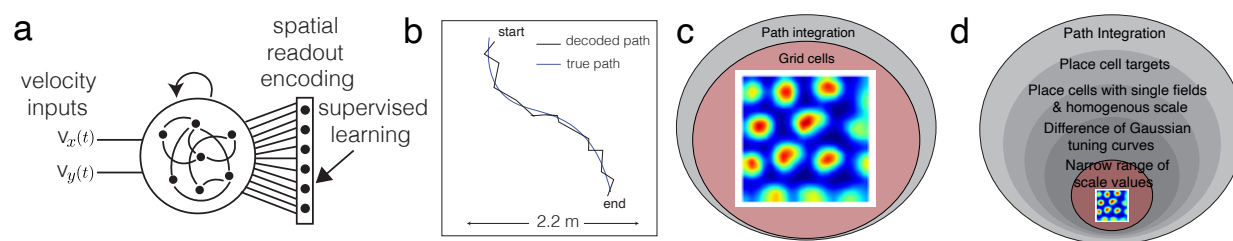


Figure 1: (a-b) Schematic of approach of training recurrent networks to predict (different possible encodings of) 2D position from 2D velocity, in a supervised manner. (b-c) Grid cells are shown to emerge in recent DL papers as a consequence of optimizing a path integration objective [12, 2, 45, 53, 36], with the suggestion that path integration implies grid cells. (d) In the current work, we show how most ANNs trained to path integrate can do so, but only a very small set of output encoding functions (vanishingly small in the function space) and very small fraction of hyperparameter space yields grid cells, leading to the conclusion that grid cell emergence results in ANNs are post hoc: they result from tuning hyperparameters to bake grid cells into networks.

novel algorithms have been shown to instead stem from seemingly minor or unstated implementation choices [51, 15, 25]. In this paper, we ask whether Neuroscientists should similarly be cautious that DL-based models of neural circuits that make specific claims about revealing the brain's optimization functions or that predict specific neural tuning curves may tell us less about fundamental scientific truths and more about programmers' particular implementation choices, and might as a result be more post hoc than predictive.

To explore these questions and point toward a systematic and circumspect use of DL for Neuroscience, we evaluate recent DL-based models of grid cells in the entorhinal-hippocampal circuit. The medial entorhinal cortex (MEC) and hippocampus (HPC) are part of the hippocampal formation, a critical brain structure for learning and memory. In a pair of Nobel-prize winning discoveries, HPC was shown to contain **place cells** [37] and MEC, its cortical input, was shown to contain **grid cells** [22]. Place cells fire at one or several seemingly random locations in small and large environments [39], respectively, while grid cells fire in a spatially periodic pattern in all two-dimensional environments, whenever the animal is at the vertices of a hexagonal lattice [22]. Over five decades, the hippocampal formation has been central to understanding how the brain organizes spatial and episodic memory, for experimentalists and theorists alike, with many mysteries remaining. A recent series of DL-based models of the circuit [12, 2, 45, 53, 36] present a story that training neural circuits on **path integration (PI)** (i.e., the task of estimating one's spatial position in an environment by integrating velocity estimates) results in the emergence of grid cells.

We use public code from prior publications to demonstrate these results are due to implementation details that tell us more about those choices than they do about MEC. By leveraging theoretically-guided large-scale hyperparameter exploration and hypothesis-driven experimentation, we show:

1. Networks trained on path integration tasks almost always learn to optimally encode position, but almost never learn grid-like representations to do so.
2. The emergence of grid-like representations depends wholly on a specifically chosen encoding of the supervised target, not on the task itself, and these choices may be unrealistic.
3. The chosen encoding requires many other sensitive hyperparameter choices, such that small alterations result in loss of grid-like representations.
4. Grid periods and period ratios depend on hyperparameter choices and are not set by the task.
5. Multiple grid modules, a fundamental characteristic of the grid cell system, do not emerge.

DL produces grid cells and some attendant properties only after making many specific design choices and searching hyperparameter space to obtain such representations, effectively baking grid cells into the task-trained networks. **It is highly improbable that DL models of path integration would have produced grid cells as a novel prediction from task-training, had grid cells not already been known to exist.** Moreover, it is unclear what interpretability or understanding these models contribute, beyond or even up to what has already been shown for these particular circuits. These results challenge the notion that deep networks offer a free lunch for Neuroscience in terms of discovering the brain's optimization

problems or generating novel a priori predictions about single-neuron representations, and warn that caution is needed when building and interpreting such models.

Our work benefited greatly from previous publications that published their code, for which the authors should be commended. By building on their code, we have been able to present novel insights that we hope will contribute to a clearer understanding of the potential risks and rewards of using and interpreting DL models in Neuroscience. To facilitate further research, we similarly release code¹.

Background: grid cells

Grid cells [22] are found in the medial entorhinal cortex of mammals and are tuned to represent the spatial location of the animal as it traverses 2D space. Each cell fires at every vertex of a triangular lattice that tiles the explored space, regardless of the speed and direction of movement through the space. As a population, grid cells exhibit several striking properties that provide support for a specialized and modular circuit. Grid cells form discrete modules (clusters), such that all cells within a module share a common period and orientation, while different modules express discretely different spatial periods [48]. The period ratios of successive modules have values in the range of 1.2-1.5.

The mechanism underlying grid cells is widely supported to be through attractor dynamics: Translation-invariant lateral connectivity within the grid cell network results in a linear Turing instability and pattern formation [7, 17, 5]. These models explain how grid cells can convert velocity inputs into updated spatial estimates, and make several predictions that have been confirmed in experiments, including most centrally the stability of low-dimensional cell-cell relationships regardless of environment and behavioral state, that define a toroidal attractor dynamics [18, 54, 50, 20, 19].

Experimental approach

Path integration (PI) is the task of using self-velocity estimates to track one’s spatial position, a crucial component of spatial navigation. The central message of existing DL models of grid cells is that training ANNs to path integrate causes the networks to learn grid cells [12, 2, 45, 53, 36].

We follow the setup used by many previous papers: a 2.2m x 2.2m arena is created, then, spatial trajectories (i.e. sequences of positions and velocities) are sampled. Networks receive as inputs the initial position and velocities, and are trained to output (a possible encoding of) the positions (Fig. 1ab) in a supervised manner. There are multiple possible encodings of position, and as we will show, this choice is critical. Two simple encodings are Cartesian [12] or polar [1]. Another encoding scheme is via bump functions in 2D space, with each output assigned a single different position that tiles the space and all outputs assigned identical tuning curves [2, 45, 45, 36]. This encoding has been equated with place cells, even though place fields tend to be heterogeneous in size and shape [39, 13], as well as in number [39], and unlike the choice of a difference-of-Gaussians (DoG) readout tuning in ANN models, do not exhibit any inhibitory surround. See Appendix 1 for position encoding details. For all encodings, supervised learning is used to train the network via backpropagation through time.

Spatial tuning assessments The spatial tuning of hidden units in the networks, visualized as ratemaps, is the primary method for comparison with the brain’s grid cells. To compute ratemaps, a trained network is evaluated on long trajectories that cover the 2D environment, then each hidden unit’s activity is binned against the true position in the environment, and counts per bin are normalized by the number of visits to that bin. Ratemaps of units are compared with grid cells through the gridness score. **We are extremely lenient with classifying a training run with particular hyperparameters a success: if even a *single* hidden unit has a grid score above a certain threshold, we say the model possesses grid cells.** The grid score is not perfect since cells classified by grid scores represent only an upper bound on the total number of grid cells; for details, see Appendices 1 and 1.

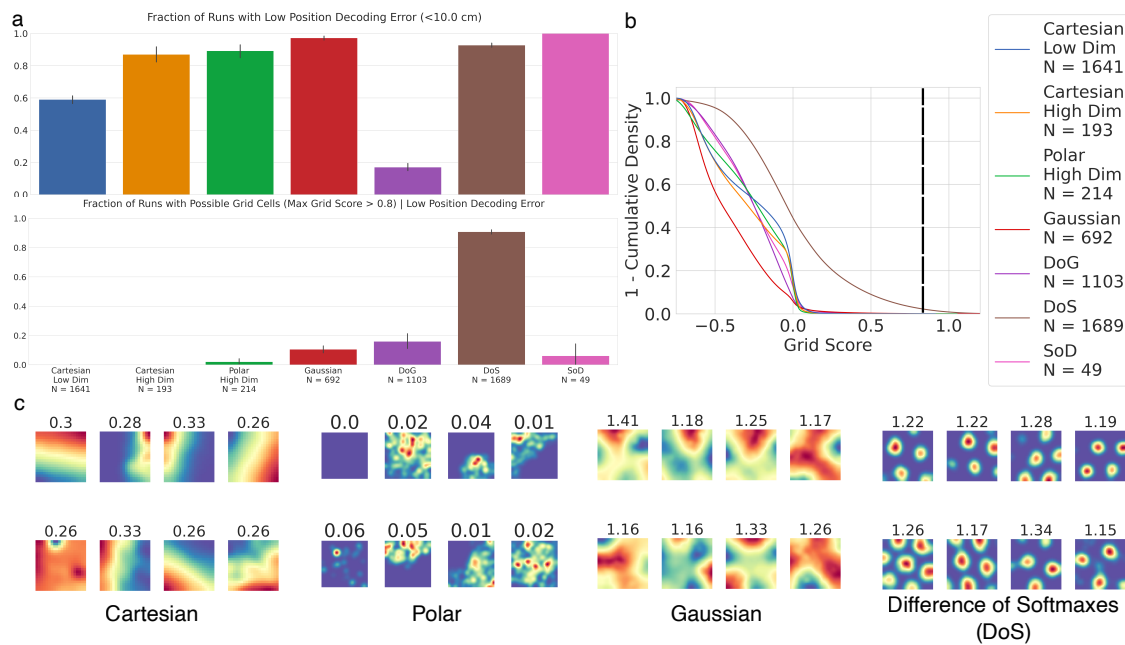


Figure 2: (a) Top: Across readout encodings, almost all networks learn to optimally encode position. Bottom: Few networks display possible grid-like representations (grid score threshold = 0.8). (b) Kernel density estimates of grid scores per readout encoding. (c) Rate maps of high grid-scoring units in deep networks trained on i) Cartesian, ii) Polar, iii) Gaussian, iv) specifically selected (tuned) Difference-of-Softmaxes (DoS) readouts. i)-iii) do not learn any grid cells. (b) Only networks trained on DoS readouts display grid-like cells. Numbers above rate maps are grid scores.

Networks trained on path integration tasks learn to estimate position, but rarely learn grid cells

As [45] wrote in their section titled “Optimally encoding position yields diverse grid-like patterns,” “Why do these diverse architectures, across diverse tasks (both navigation and autoencoding), all converge to a grid-like solution, and what governs the lattice structure of this solution?” We demonstrate, in contrast, that most networks do not converge to a grid-like solution, instead requiring very specific readout tuning functions and hyperparameter choices. Grid-like representation emergence is determined by choices made by the programmer that, rather than relating to the path integration objective or biologically realistic place cells, is designed post hoc to produce grid cells.

We ran large-scale hyperparameter sweeps across common implementation choices: 1) Architectures: RNN [14]; LSTM [24]; GRU [10]; UGRNN [11]; 2) Activation: Sigmoid; Tanh; ReLU; Linear; 3) Optimizers: SGD, Adam [31]; RMSProp [23] 4) Supervised Targets: Cartesian; Polar; high-dimensional bump-like readout population code with Gaussian [2], **Difference-of-Softmaxes (DoS)** [45, 46, 36] or Difference-of-Gaussians (DoG) tuning curves. 5) Loss: mean squared error on the agent’s Cartesian position [27, 12]; geodesic distance on the agent’s polar position [1]; cross entropy on a high-dimensional population of bump-like readout units [2, 45, 36] 6) Miscellaneous: recurrent & readout dropout, initialization, regularization, seed.

For networks with bump-like population readouts, we additionally swept: 1) Field width σ i.e. standard deviation of the Gaussian tuning curve (often denoted $\sigma_E \stackrel{\text{def}}{=} \sigma$ in the literature); 2) Whether the bump-like readouts have homogeneous or heterogeneous field widths; 3) In the case of readouts with DoG or DoS tuning curves, the surround scale s i.e. the ratio between the inhibitory and excitatory Gaussian’s standard deviations ($s \stackrel{\text{def}}{=} \sigma_I / \sigma_E$); 4) Number of fields per readout unit. Evaluating the entire hyperparameter volume is computationally prohibitive, so we evaluated a subvolume most consistent with previous papers, focusing on perturbations from hyperparameters that did produce grid cells. In this sense, **our search was biased toward configurations shown to produce grid cell emergence and thus our findings about the fragility of these solutions conservatively favored these solutions as much as possible.** All sweeps

¹<https://github.com/RylanSchaeffer/MEC-HPC-Models-Investigation>

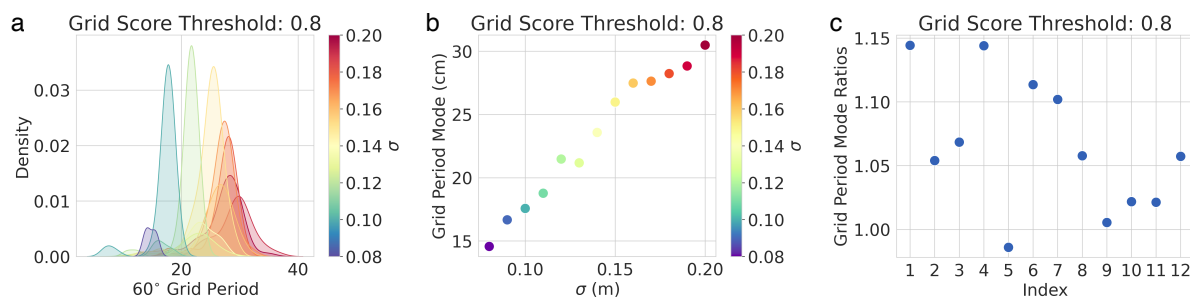


Figure 3: Grid periods are hyperparameter-dependent and multiple modules do not emerge. (a) Over a wide sweep of DoS (the most likely to produce grid cells, Fig. 2a bottom) target field widths σ_E , the distribution of grid periods is unimodal (each color is distribution of periods obtained from 3 runs with a σ_E value; different periods are only obtained by varying σ_E), meaning multiple grid modules do not emerge, in contrast to real grid cells. (b) The chosen target field width σ_E determines the grid period mode, meaning that hand-designed hyperparameter choices, not an intrinsic emergent property, sets the grid period. (c) If we use grid periods obtained from smoothly sweeping σ_E as a proxy for different modules, the period ratios are closer to 1 than the experimental ratios of ~ 1.4 .

are provided in Appendix 1.

To evaluate whether a network learns to optimally encode its position, we measured its position decoding error using previous papers' method [2, 45, 36]: we computed position decoding error using the network's output Cartesian positions (if trained on Cartesian targets) or by decoding position from the network's outputs. Any network with error < 10 cm was considered to have achieved optimal position encoding; this threshold was chosen based on noise inherent in the decoding algorithm.

In total, we trained > 5500 networks and found that most hyperparameter configurations succeed in learning to path integrate (Fig. 2a, Top), but few learn grid cell representations (Fig. 2a, Bottom). This is consistent with earlier work [27, 1] demonstrating that networks can learn to path integrate and solve other navigational problems (e.g. estimating which of several environments correspond to the current location) without lattice cells emerging as a solution.

Lattice cell emergence requires a highly specific choice of supervised target encoding

We next sought to characterize when grid cells are learnt under different encodings of the supervised target i.e., 2D position. We tested multiple encodings: i) Cartesian, ii) Polar, iii) Gaussian, iv) a very specifically selected, Difference-of-Gaussian (DoG/DoS) readout and more. We found that grid cells do not emerge from Cartesian, Polar or Gaussian encodings (Fig. 2abc), consistent with earlier work [27]. Only by choosing a DoG readout tuning curve did grid cells sometimes emerge.

Grid periods are parameter-dependent and multiple modules do not emerge

Next, a prominent feature of grid cells that is critical for their unambiguous encoding of position over large scales is the existence of a discrete set of grid periodicities, which tend to scale by a rough factor of 1.4 between adjacent scales [48]. We asked whether ANN models generate multiple periods and when they did so, what hyperparameters and other choices the formed periods depended on.

To ensure we would obtain at least some grid cells, we fixed the readouts to be DoSs, and swept over different scales of the DoS. We found that almost all runs had a unimodal distribution of grid periods (Fig. 3a), meaning the networks learnt only one module of grid cells.

Further, we found that the period of the formed grid-like representation is completely determined by the scale σ_E of the externally imposed readout DoG (Fig. 3). The period of the grid-like responses in every run increased monotonically with the width of the DoG readout (Fig. 3b). Since the models did not result in multiple modules in a single network, we used the somewhat discrete distribution of peaks of the single module formed when sweeping the DoG parameter more continuously to compute grid period ratios. These period ratios from adjacent peaks led to non-biological values (Fig. 3c).

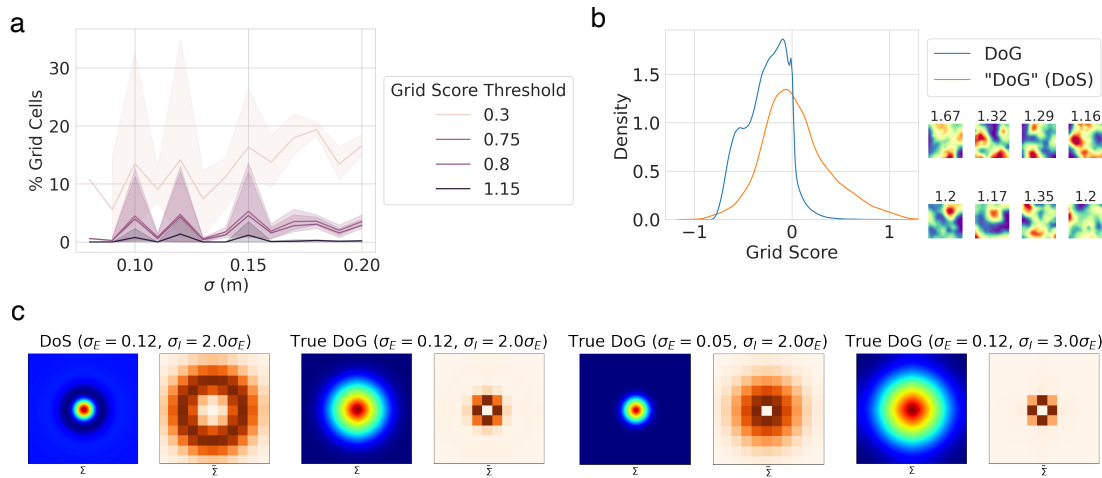


Figure 4: Other changes to Difference-of-Softmaxes (DoS) affect the formation of grid cells: (a) The existence of grid solutions even with DoG readout encodings is highly sensitive to parameters such as the target function receptive field width σ_E , and small alterations result in the disappearance of grid cells regardless of grid score threshold. (b) (left) Comparison of grid scores of trained networks with DoG shows that the Difference-of-Softmaxes (DoS), as used by [45, 36, 46], is critical for high grid scores, (right) Rate maps of highest scoring units from DoG networks achieving low position decoding error. (c) Computing the Fourier transform (right of each pair) of the readout second-moment matrix (left of each pair) explains that DoS places Fourier power on an annulus with sufficiently large radius to produce grid cells, showing why the particular DoS choice is critical.

Fourier analysis of Turing instability explains the preceding empirical results Why do only Difference-of-Softmaxes (DoS) readouts produce grid cells? We restate the essence of previous analyses of first-principles models [7, 29] here, and explain why they might explain our empirical findings.

In a recurrent network with dynamics $\dot{r}(x) = -r(x) + g(W \star r)$, where x designates the neural index (in a continuum approximation for neurons), $W \star r$ designates the total (integrated) inputs from the network to the neuron at index x , and g is the non-linearity, if the recurrent weight interaction is translationally invariant, then $W(x, x') = W(x - x') = W(\Delta x)$. Under DoG interactions:

$$W(\Delta x) \equiv f(\Delta x) = \alpha_E \exp\left(-\frac{(\Delta x)^2}{2\sigma_E^2}\right) - \alpha_I \exp\left(-\frac{(\Delta x)^2}{2\sigma_I^2}\right) \quad (1)$$

where Δx refers to the difference of indices between the neural pair linked by the weights. The evolution of activity can be decomposed into the growth and decay of Fourier components of the rate vector, which is fully determined by the Fourier transform of W , which is given by:

$$\tilde{f}(k) = \int_{\mathcal{R}} d(\Delta x) f(\Delta x) e^{ik\Delta x} = \alpha_E \sigma_E \exp\left(-\frac{\sigma_E^2 k^2}{2}\right) - \alpha_I \sigma_I \exp\left(-\frac{\sigma_I^2 k^2}{2}\right) \quad (2)$$

Here α_E (α_I) denotes the strength and σ_E (σ_I) denotes the scale of excitation (inhibition). For linearized dynamics that approximate $\dot{r}(x) \sim -r(x) + f(\Delta x) \star r$ (i.e. g has been linearized), the solution will be periodic if the maxima of $\tilde{f}(k)$, given by $[k^*]^2 = \frac{2}{\sigma_E^2 - \sigma_I^2} \log(\alpha_E \sigma_E^3 / \alpha_I \sigma_I^3)$, contains sufficient power and if $k^* \neq 0$. Specifically, the condition for pattern formation is $\tilde{f}(k^*) > 1$ [8, 29]. In particular, the inhibitory surround contained in $f(\Delta x)$, with strength σ_I , is key to pattern formation; if $\sigma_I \rightarrow \infty$ or $\alpha_I \rightarrow 0$, the maximum is at the origin ($k^* = 0$), causing no pattern formation. Therefore, a Gaussian interaction cannot produce periodic patterns.

The connection of this theory to ANNs trained with supervised readout target P was made in [45] through the observation that dynamics with an MSE loss $\|P - W_{\text{readout}} r\|^2$ can be approximated as $\dot{r} = -\lambda r + \Sigma r$, where $\Sigma = PP^T$ is the readout correlation matrix and λ is a regularization parameter. This dynamics now looks identical to the first dynamics with recurrent interaction matrix $W(x, x') = \Sigma_{xx'}$. For identical, unimodal target readout functions with DoG tuning curves, this readout correlation matrix is also a difference

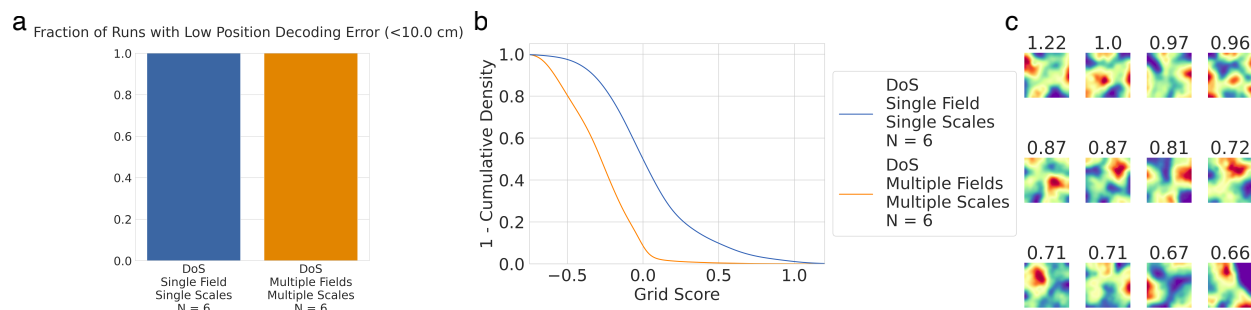


Figure 5: Adding place cell-like heterogeneity to the readouts prevents the formation of grid cells. To maximally favor the emergence of grid representations, we selected DoS encodings with the best hyperparameters for grid cell emergence (RNN or UGRNN, ReLU, $\sigma_E = 0.12$ cm, $s = 2.0$, 3 seeds), then tested the effect of multiple fields per place cell ($\sim 1 + \text{Pois}(3.0)$) and multiple scales (receptive field width $\sigma_E \sim \text{Unif}(0.06, 1.0)$ m and surround scale $s \sim \text{Unif}(1.25, 4.5)$). (a) Networks trained on single-field single-scale DoS and multi-scale multi-field DoS (more similar to biological place cells) readouts all obtain low position decoding error. (b) Multi-scale multi-field DoS readouts do not learn grid cells. (c) Highest-scoring rate maps from multi-field multi-scale networks.

of Gaussians (Appendix 1). By the same theory outlined above, Gaussian tuning curves should not produce periodic patterns.

Unmentioned implementation details matter We discovered a seemingly minor implementation detail in several preceding papers [45, 46, 36] that proved important for the emergence of grid cells, that (to the best of our knowledge) is unmentioned in the main texts and supplements. While these papers nominally refer to a Difference-of-Gaussian (DoG) readout target function (Eqn. 1), their code instead uses a Difference-of-Softmaxes (DoS) target function. When we trained ideal grid-forming ReLU networks on place-cell-like target encodings with DoG tuning curves, sweeping the receptive field σ and surround scale s , DoG readouts did not result in grid cells (Fig. 4b), while DoS readouts did. The preceding Fourier analysis shows why DoS tuning, though still more unrealistic as a place cell model than DoG, matters (Fig. 4c).

Grid cells disappear with more biologically-realistic readout functions

Place cells, to which grid cells project, differ significantly from the highly simplified single-scale, single-field ANN readouts. Place cells have heterogeneously distributed field widths, many with multiple fields [39, 13] (and are not DoS-like functions). This naturally leads to the question: Will readouts targets with multiple scales and multiple fields per place cell still produce grid cells?

We found that networks trained with multiple-field multiple-scale bump encodings per readout unit achieve comparably low position decoding error as good single-field single-scale encodings (Fig. 5a), but do not learn grid cells (Fig. 5bc). This further demonstrates that accurate position encoding and integration does not require grid representations.

Discussion

For research that uses deep networks as models of the brain, there is a fundamental obstacle to making the claim that a given optimization problem is what the brain is solving: If we know the responses of a significant fraction of units from biological networks performing a certain task, we cannot infer the loss function that the brain is optimizing since in principle, numerous different loss functions can have the same/similar minima (Fig. 6 top). In other words, there is typically a many-to-one mapping between loss functions and some point in state space where the functions have a minimum. Some of the different grid models from deep learning and first principles show that this is possible [12, 2, 45, 53, 7]. Conversely, given a reasonable optimization problem that we select based on an organism's ecological niche, we cannot infer a single solution (and thus build truly predictive single-cell tuning models), since there exist several minima to that optimization problem (Fig. 6 bottom). In other words, there is typically a one-to-many mapping from a loss function to its set of

solutions. To break this degeneracy of multiple minima and arrive at truly predictive models or a better understanding of the brain's optimization problems, we must acknowledge, understand, and model the specific inductive biases and constraints present in the biological system we are trying to model. It is untenable to expect success without doing so. This is what we refer to as an informal neural 'no-free-lunch'. In these cases, these constraints and biases are as informative as other design principles, such as the loss function.

What can we learn from DL models about brain circuits without considering and studying biological inductive biases? Population-wide low-dimensional latent representations and dynamics that arise as necessary for solving difficult problems are possibly sufficiently robust and abstracted to be predictive of population dynamics in a neural circuit without addition of detailed biological constraints. This can explain successes of the population-level analyses of the visual pathway [28, 3], as well as the population-level low-dimensional dynamics of circuits solving inference tasks [49, 42, 1, 52]. In these cases, the emergent solutions are fundamental features of any systems that solve the task, and by construction need not be specific to brains.

Returning to grid cells, we have conclusively shown that they do not generically arise in networks trained to path-integrate. This raises the question of what different additional architectural, hyperparameter, and constraint choices had to be made in previous papers [2, 45, 46] to obtain grid-like tuning, given that they used a path integration objective with Gaussian place cell targets, requirements that are reasonable but not sufficient to obtain grid-like tuning.

Theoretical work on grid cell representations [16, 47, 34] suggests that the following two factors are important features of the grid cell code: a very large coding range and the related property of robustness/intrinsic error correcting coding. Adding these properties to the loss are likely to be a more principled way to obtain grid-like fields rather than by hand-designing a biologically unrealistic DoG or DoS readout. Consistent with this, the very large amount of dropout or other noise required in [2, 12] suggests that the coding-theoretic insight on intrinsic error-correction properties of grid cells and their related large capacity may indeed be central ingredients to produce grid cells.

In fact, there are a number of key properties of the grid cell code elucidated earlier by theoretical arguments, including but not limited to path integration, that we hypothesize may form a sufficient and biologically relevant set for the emergence of grid cells: 1) non-negative activations; 2) equivariant population responses (i.e., the population response always lies on a hypersphere); 3) a path integrating (PI) code or in other words, translation invariant representations [17, 7]; 4) high representational capacity [6, 47, 33]; 5) intrinsic error-correcting capabilities [6, 47]; and finally, 6) uniformly distributed (whitened) and low spatial information per cell, such that the total spatial information of the grid code should be equally distributed across all cells.

We suggest that several of these properties (e.g. excluding the more specific property 3)) are extremely general features of neural codes, and incorporating them into ANN models of neural circuits could increase their ability to make de novo rather than post hoc predictions, and remove the need to use synthetic and unrealistic choices with fine hyperparameter tuning (such as DOG and DOS readouts).

In sum, the findings here argue that the central contribution of scientific interest when building and studying ANN models of the brain that reproduce specific tuning curves is not that the model produced the curves (after all, given the expressive power of deep networks, it is no surprise that training them to effectively generate a given tuning will in fact succeed), but rather to explore and carefully characterize conditions and constraints under which the particular tuning does and does not emerge, to consider whether the constraints align with biological constraints, and to consider what principles can be extracted from them. When not predicting specific tuning curves, but rather population-level response dynamics, the results are expected to be more robust to specific implementational choices and thus more durable.

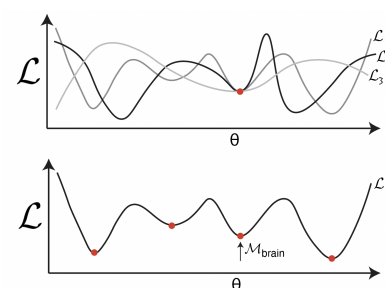


Figure 6: Challenges in achieving the two central claims of recent DL models of neuroscience: Top: Building a model that replicates observed neural responses does not guarantee that the loss function used is the brain's objective, as multiple objectives can share a solution. Bottom: Training a network on a plausible loss function or even the correct loss function need not yield the solution the brain has selected because the loss function may have multiple minima, of which the brain selects one based on its constraints, while an ANN selects another, based on the optimization technique used.

References

- [1] Emergence of dynamically reconfigurable hippocampal responses by learning to perform probabilistic spatial reasoning. *bioRxiv*.
- [2] Andrea Banino, Caswell Barry, Benigno Uribe, Charles Blundell, Timothy Lillicrap, Piotr Mirowski, Alexander Pritzel, Martin J. Chadwick, Thomas Degris, Joseph Modayil, Greg Wayne, Hubert Soyer, Fabio Viola, Brian Zhang, Ross Goroshin, Neil Rabinowitz, Razvan Pascanu, Charlie Beattie, Stig Petersen, Amir Sadik, Stephen Gaffney, Helen King, Koray Kavukcuoglu, Demis Hassabis, Raia Hadsell, and Dharshan Kumaran. Vector-based navigation using grid-like representations in artificial agents. *Nature*, 557(7705):429–433, May 2018.
- [3] Pouya Bashivan, Kohitij Kar, and James J. DiCarlo. Neural population control via deep image synthesis. *Science*, 364(6439):eaav9436, May 2019. Publisher: American Association for the Advancement of Science.
- [4] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language Models are Few-Shot Learners. *arXiv:2005.14165 [cs]*, July 2020. arXiv: 2005.14165.
- [5] Yoram Burak and Ila Fiete. Do We Understand the Emergent Dynamics of Grid Cell Activity? *Journal of Neuroscience*, 26(37):9352–9354, September 2006. Publisher: Society for Neuroscience Section: Journal Club.
- [6] Yoram Burak and Ila R Fiete. Unpublished observations. 2008.
- [7] Yoram Burak and Ila R. Fiete. Accurate Path Integration in Continuous Attractor Network Models of Grid Cells. *PLOS Computational Biology*, 5(2):e1000291, February 2009. Publisher: Public Library of Science.
- [8] Yoram Burak and Ila R Fiete. Accurate path integration in continuous attractor network models of grid cells. *PLoS Comput Biol*, 5(2):e1000291, Feb 2009.
- [9] Malcolm G. Campbell, Samuel A. Ocko, Caitlin S. Mallory, Isabel I. C. Low, Surya Ganguli, and Lisa M. Giocomo. Principles governing the integration of landmark and self-motion cues in entorhinal cortical codes for navigation. *Nature Neuroscience*, 21(8):1096–1106, August 2018. Number: 8 Publisher: Nature Publishing Group.
- [10] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling. *NIPS Workshop Deep Learning and Representation Learning*, December 2014. Number: arXiv:1412.3555 arXiv:1412.3555 [cs].
- [11] Jasmine Collins, Jascha Sohl-Dickstein, and David Sussillo. Capacity and Trainability in Recurrent Neural Networks. *International Conference on Learning Representations*, March 2017. Number: arXiv:1611.09913 arXiv:1611.09913 [cs, stat].
- [12] Christopher J Cueva and Xue-Xin Wei. Emergence of grid-like representations by training recurrent neural networks to perform spatial localization. *International Conference on Learning Representations*, page 19, 2018.
- [13] Tamir Eliav, Shir R. Maimon, Johnatan Aljadeff, Misha Tsodyks, Gily Ginosar, Liora Las, and Nachum Ulanovsky. Multiscale representation of very large environments in the hippocampus of flying bats. *Science*, 372(6545):eabg4020, May 2021. Publisher: American Association for the Advancement of Science.
- [14] Jeffrey L. Elman. Finding Structure in Time. *Cognitive Science*, 14(2):179–211, 1990. _eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1207/s15516709cog1402_1.
- [15] Logan Engstrom, Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Firdaus Janoos, Larry Rudolph, and Aleksander Madry. Implementation Matters in Deep Policy Gradients: A Case Study on PPO and TRPO. *arXiv:2005.12729 [cs, stat]*, May 2020. arXiv: 2005.12729.

- [16] Ila R Fiete, Yoram Burak, and Ted Brookings. What grid cells convey about rat location. *J Neurosci*, 28(27):6858–71, Jul 2008.
- [17] Mark C Fuhs and David S Touretzky. A spin glass model of path integration in rat medial entorhinal cortex. *J Neurosci*, 26(16):4266–4276, 2006.
- [18] Marianne Fyhn, Torkel Hafting, Alessandro Treves, May-Britt Moser, and Edvard I Moser. Hippocampal remapping and grid realignment in entorhinal cortex. *Nature*, 446(7132):190–194, 2007.
- [19] Richard J. Gardner, Erik Hermansen, Marius Pachitariu, Yoram Burak, Nils A. Baas, Benjamin A. Dunn, May-Britt Moser, and Edvard I. Moser. Toroidal topology of population activity in grid cells. *Nature*, 602(7895):123–128, February 2022. Number: 7895 Publisher: Nature Publishing Group.
- [20] Richard J. Gardner, Li Lu, Tanja Wernle, May-Britt Moser, and Edvard I. Moser. Correlation structure of grid cells is preserved during sleep. *Nature Neuroscience*, 22(4):598–608, April 2019. Number: 4 Publisher: Nature Publishing Group.
- [21] Joshua I. Glaser, Ari S. Benjamin, Raed H. Chowdhury, Matthew G. Perich, Lee E. Miller, and Konrad P. Kording. Machine Learning for Neural Decoding. *eNeuro*, 7(4), July 2020. Publisher: Society for Neuroscience Section: Research Article: Methods/New Tools.
- [22] Torkel Hafting, Marianne Fyhn, Sturla Molden, May-Britt Moser, and Edvard I. Moser. Microstructure of a spatial map in the entorhinal cortex. *Nature*, 436(7052):801–806, August 2005. Number: 7052 Publisher: Nature Publishing Group.
- [23] Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. Lecture 6e-RMSProp.
- [24] Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9(8):1735–1780, November 1997.
- [25] Andrew Ilyas, Logan Engstrom, Shibani Santurkar, Dimitris Tsipras, Firdaus Janoos, Larry Rudolph, and Aleksander Madry. A Closer Look at Deep Policy Gradients. *arXiv:1811.02553 [cs, stat]*, May 2020. arXiv: 1811.02553.
- [26] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislaw Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steinegger, Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstein, David Silver, Oriol Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate protein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, August 2021. Number: 7873 Publisher: Nature Publishing Group.
- [27] I. Kanitscheider and I. R. Fiete. Training recurrent networks to generate hypotheses about how the brain solves hard navigation problems. *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [28] Alexander J. E. Kell, Daniel L. K. Yamins, Erica N. Shook, Sam V. Norman-Haignere, and Josh H. McDermott. A Task-Optimized Neural Network Replicates Human Auditory Behavior, Predicts Brain Responses, and Reveals a Cortical Processing Hierarchy. *Neuron*, 98(3):630–644.e16, May 2018.
- [29] Mikail Khona, Sarthak Chandra, and Ila R. Fiete. From smooth cortical gradients to discrete modules: spontaneous and topologically robust emergence of modularity in grid cells. *bioRxiv*, page 2021.10.28.466284, January 2022.
- [30] Timothy D. Kim, Thomas Z. Luo, Jonathan W. Pillow, and Carlos D. Brody. Inferring Latent Dynamics Underlying Neural Population Activity via Neural Differential Equations. In *Proceedings of the 38th International Conference on Machine Learning*, pages 5551–5561. PMLR, July 2021. ISSN: 2640-3498.
- [31] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *International Conference on Learning Representations*, January 2017. Number: arXiv:1412.6980 arXiv:1412.6980 [cs].

- [32] Jesse A. Livezey, Kristofer E. Bouchard, and Edward F. Chang. Deep learning as a tool for neural data analysis: Speech classification and cross-frequency coupling in human sensorimotor cortex. *PLOS Computational Biology*, 15(9):e1007091, September 2019. Publisher: Public Library of Science.
- [33] A. Mathis, A. Herz, and M. Stemmler. Optimal population codes for space: grid cells outperform place cells. *Neural Comp.*, 24:2280–2317, 2012.
- [34] Alexander Mathis, Andreas V. M. Herz, and Martin Stemmler. Optimal Population Codes for Space: Grid Cells Outperform Place Cells. *Neural Computation*, 24(9):2280–2317, September 2012.
- [35] Alexander Mathis, Pranav Mamidanna, Kevin M. Cury, Taiga Abe, Venkatesh N. Murthy, Mackenzie Weygandt Mathis, and Matthias Bethge. DeepLabCut: markerless pose estimation of user-defined body parts with deep learning. *Nature Neuroscience*, 21(9):1281–1289, September 2018. Number: 9 Publisher: Nature Publishing Group.
- [36] Aran Nayeibi, 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 *Advances in Neural Information Processing Systems*, volume 34, pages 12167–12179. Curran Associates, Inc., 2021.
- [37] J. O’Keefe and J. Dostrovsky. The hippocampus as a spatial map: Preliminary evidence from unit activity in the freely-moving rat. *Brain Research*, 34:171–175, 1971. Place: Netherlands Publisher: Elsevier Science.
- [38] Talmo D. Pereira, Diego E. Aldarondo, Lindsay Willmore, Mikhail Kislin, Samuel S.-H. Wang, Mala Murthy, and Joshua W. Shaevitz. Fast animal pose estimation using deep neural networks. *Nature Methods*, 16(1):117–125, January 2019. Number: 1 Publisher: Nature Publishing Group.
- [39] P. D. Rich, H.-P. Liaw, and A. K. Lee. Large environments reveal the statistical structure governing hippocampal representations. *Science*, 345(6198):814–817, August 2014.
- [40] Muhammad Saif-ur Rehman, Robin Lienkämper, Yaroslav Parpaley, Jörg Wellmer, Charles Liu, Brian Lee, Spencer Kellis, Richard Andersen, Ioannis Iossifidis, Tobias Glasmachers, and Christian Klaes. SpikeDeeptector: a deep-learning based method for detection of neural spiking activity. *Journal of Neural Engineering*, 16(5):056003, July 2019. Publisher: IOP Publishing.
- [41] Andrew Saxe, Stephanie Nelli, and Christopher Summerfield. If deep learning is the answer, what is the question? *Nature Reviews Neuroscience*, 22(1):55–67, January 2021. Number: 1 Publisher: Nature Publishing Group.
- [42] Rylan Schaeffer, Mikail Khona, Leenoy Meshulam, Brain Laboratory International, and Ila Fiete. Reverse-engineering recurrent neural network solutions to a hierarchical inference task for mice. *Advances in Neural Information Processing Systems*, 33:4584–4596, 2020.
- [43] David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, January 2016. Number: 7587 Publisher: Nature Publishing Group.
- [44] Trygve Solstad, Charlotte N. Boccara, Emilio Kropff, May-Britt Moser, and Edvard I. Moser. Representation of Geometric Borders in the Entorhinal Cortex. *Science*, 322(5909):1865–1868, December 2008. Publisher: American Association for the Advancement of Science.
- [45] Ben Sorscher, Gabriel C Mel, Surya Ganguli, and Samuel A Ocko. A unified theory for the origin of grid cells through the lens of pattern formation. *Advances in Neural Information Processing Systems*, page 18, 2019.
- [46] Ben Sorscher, Gabriel C. Mel, Samuel A. Ocko, Lisa Giocomo, and Surya Ganguli. A unified theory for the computational and mechanistic origins of grid cells. Technical report, bioRxiv, December 2020. Section: New Results Type: article.

- 429 [47] Sameet Sreenivasan and Ila Fiete. Grid cells generate an analog error-correcting code for singularly
430 precise neural computation. *Nat Neurosci*, 14(10):1330–7, Sep 2011.
- 431 [48] Hanne Stensola, Tor Stensola, Trygve Solstad, Kristian Frøland, May-Britt Moser, and Edvard I. Moser.
432 The entorhinal grid map is discretized. *Nature*, 492(7427):72–78, December 2012. Number: 7427 Publisher:
433 Nature Publishing Group.
- 434 [49] David Sussillo and Omri Barak. Opening the Black Box: Low-Dimensional Dynamics in High-Dimensional
435 Recurrent Neural Networks. *Neural Computation*, 25(3):626–649, March 2013.
- 436 [50] S.G. Trettel, J.B. Trimper, E. Hwaun, I.R. Fiete, and L.L. Colgin. Grid cell co-activity patterns during
437 sleep reflect spatial overlap of grid fields during active behaviors. *Nat Neurosci*, 22(4):609–617, 04 2019.
- 438 [51] George Tucker, Surya Bhupatiraju, Shixiang Gu, Richard E. Turner, Zoubin Ghahramani, and Sergey
439 Levine. The Mirage of Action-Dependent Baselines in Reinforcement Learning. *arXiv:1802.10031 [cs,*
440 *stat]*, November 2018. arXiv: 1802.10031.
- 441 [52] Jakob Voigts, Ingmar Kanitscheider, Nicholas J. Miller, Enrique H. S. Toloza, Jonathan P. Newman,
442 Ila R. Fiete, and Mark T. Harnett. Spatial reasoning via recurrent neural dynamics in mouse retrosplenial
443 cortex. *bioRxiv*, April 2022.
- 444 [53] James C. R. Whittington, Timothy H. Muller, Shirley Mark, Guifen Chen, Caswell Barry, Neil Burgess,
445 and Timothy E. J. Behrens. The Tolman-Eichenbaum Machine: Unifying Space and Relational Memory
446 through Generalization in the Hippocampal Formation. *Cell*, 183(5):1249–1263.e23, November 2020.
- 447 [54] K.J. Yoon, M.A. Buice, R. Barry, C. and Hayman, N. Burgess, and I.R. Fiete. Specific evidence of
448 low-dimensional continuous attractor dynamics in grid cells. *Nat Neurosci*, 16(8):1077–84, Aug 2013.

1 Checklist

1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?
- (b) Did you describe the limitations of your work?
- (c) Did you discuss any potential negative societal impacts of your work? We do not feel our paper has potential negative societal impacts.
- (d) Have you read the ethics review guidelines and ensured that your paper conforms to them?

2. If you are including theoretical results...

- (a) Did you state the full set of assumptions of all theoretical results?
- (b) Did you include complete proofs of all theoretical results?

3. If you ran experiments...

- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)?
- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)?
- (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)?
- (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? We haven't yet had time to collect this information, but we will add this information to the final version if accepted.

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- (a) If your work uses existing assets, did you cite the creators?
- (b) Did you mention the license of the assets?
- (c) Did you include any new assets either in the supplemental material or as a URL?
- (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating?
- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content?

5. If you used crowdsourcing or conducted research with human subjects...

- (a) Did you include the full text of instructions given to participants and screenshots, if applicable?
- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable?
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation?

Position Encodings

Suppose we sample a sequence of positions $x_0, \dots, x_T \in \mathbb{R}^2$ and a sequence of velocities $v_1, \dots, v_T \in \mathbb{R}^2$, where $x_t = x_{t-1} + v_t$. We want to train the networks in a supervised manner to predict (a possible encoding of) their position. We used the below encodings as different regression targets. For some encodings that required place cell populations, we denote the number of place cells N_p and denote their locations $\{p_i\}_{i=1}^{N_p}$, sampled uniformly at random within the $2.2 \text{ m} \times 2.2 \text{ m}$ arena.

- **“Low-dimensional” Cartesian:** x_t

- **“High-dimensional” Cartesian:** Let $A \in \mathbb{R}^{N_p \times 2}, b \in \mathbb{R}^{N_p}$ have entries sampled i.i.d. from $Uniform(-1, 1)$. The target is a vector in \mathbb{R}^{N_p} given by:

$$Ax_t + b$$

- **Polar:** Let (r_t, θ_t) denote the polar form of the agent’s position x_t . The target is a vector in \mathbb{R}^{N_p} , half comprised of “distance” units and half comprised of “direction” units. Let $A \in \mathbb{R}^{0.5N_p \times 1}, b \in \mathbb{R}^{0.5N_p}$ have entries sampled i.i.d. from $Uniform(-1, 1)$; the distance cells have activities:

$$Ar_t + b$$

Let $\mu \in [-\pi, \pi]^{N_p/2}$ have entries sampled i.i.d. uniformly at random. The direction cells have entries given by von-Mises-like bumps:

$$\exp\left(\cos(\theta_t - \mu_i)\right)$$

- **Gaussian:** A vector in \mathbb{R}^{N_p} whose entries are given by:

$$\exp\left(-\frac{1}{2\sigma_E^2}\|x_t - p_i\|^2\right)$$

- **Difference of Gaussians:** A vector in \mathbb{R}^{N_p} whose entries are given by:

$$\exp\left(-\frac{1}{2\sigma_E^2}\|x_t - p_i\|^2\right) - \exp\left(-\frac{1}{2\sigma_I^2}\|x_t - p_i\|^2\right)$$

- **Difference of Softmaxes:** A vector in \mathbb{R}^{N_p} whose entries are given by:

$$\text{Softmax}\left(-\frac{1}{2\sigma_E^2}\|x_t - p_i\|^2\right) - \text{Softmax}\left(-\frac{1}{2\sigma_I^2}\|x_t - p_i\|^2\right)$$

- **Softmax of Differences:** A vector in \mathbb{R}^{N_p} whose entries are given by:

$$\text{Softmax}\left(-\frac{1}{2\sigma_E^2}\|x_t - p_i\|^2 + \frac{1}{2\sigma_I^2}\|x_t - p_i\|^2\right)$$

Grid Scores and Grid Cell Thresholds

What qualifies as a grid cell? The most commonly used method of quantifying grid cells is via the “grid score”, which functions by binning neural activity into rate maps using spatial position, applying an adaptive smoother, then taking a circular sample of the autocorrelation centered on the central peak and comparing it to rotated versions of the same circular sample. The 60° grid score is specifically given by:

$$(corr[60] + corr[120])/2 - (corr[30] + corr[90] + corr[150])/3$$

We used the grid scorer implementation used by [2] (<https://github.com/deepmind/grid-cells/blob/master/scores.py>), [45] (<https://github.com/ganguli-lab/grid-pattern-formation/blob/master/scores.py>) and [36].

What score is sufficient to qualify as a grid score? Experimentalists have used thresholds of 0.3 [44] and 0.349 [9] on biological neurons, whereas computationalists have used 0.3 [36] and 0.37 [2, 45] on artificial neurons. We found that for artificial neurons, these thresholds are far too low (Fig. 7); ANN units with grid scores > 0.4 , and even as high as 1.3, often look nothing like grid cells. This is because the grid score looks for 60° rotational symmetry, and while grid cells are indeed symmetric, so are many other rate maps. Private conversations with the authors of [2, 45, 36] showed everyone was in agreement that grid scores are a metric in need of improvement. After internal disagreement, we compromised at a grid score threshold of 0.8.

Number of Bins for Computing Rate Maps

The first step in computing grid scores is determining the number of bins to use to compute rate maps. However, establishing the number of bins used by previous approaches proved baffling difficult. The original experimental work used 5 cm by 5 cm bins [22]. Since the square arena used in these experiments is 2.2 m by 2.2 m (same as [2, 45, 36], the number of bins should be 44 x 44. However, we found discrepancies in previous approaches' text and code.

[2]'s text claims to use 32 x 32 bins of 6.875 cm x 6.875 cm ("Spatial (ratemaps) and directional activity maps were calculated for individual units as follows. Each point in the trajectory was assigned to a specific spatial and directional bin according to its location and the direction in which it faced. Spatial bins were defined as a 32 x 32 square grid spanning each environment and directional bins as 20 equal width intervals."), but their code² used 20 x 20 bins of 11 cm x 11 cm. [45] claimed to use 2 cm x 2 cm bins ("Grid score was evaluated as in [2]). A spatial ratemap was computed for each neuron by binning the agent's position into 2cm x 2cm bins, and computing the average firing rate within each bin.") but their code³ similarly used 20 x 20 bins. [36] claimed to use 5 cm x 5 cm bins ("Nayebi: "We bin the positions in each environment using 5 cm bins, following prior work [Hardcastle et al., 2017, Butler et al., 2019, Low et al., 2020]. Thus, the 100cm² environment used 400 (20 x 20) bins, the 150cm² environment used 900 (30 x 30) bins, and the 400cm 1D

²<https://github.com/deepmind/grid-cells/blob/master/train.py#L201>

³<https://github.com/ganguli-lab/grid-pattern-formation/blob/master/visualize.py#L136>



Figure 7: **Grid scoring is the dominant method to identify grid cells, but we found it performs extremely poorly at identifying grid cells.** Top: example rate maps from low position decoding error Difference-of-Softmaxes (DoS) networks three grid score ranges: [0.35, 0.45], [0.8, 0.9], [1.15, ∞). Bottom: example rate maps from low position decoding error Difference-of-Gaussians (DoG) networks. We considered three grid score thresholds: 0.3 (used by some experimentalists), 0.8 (low probability of finding grid cells), 1.15 (decent probability of finding grid cells).

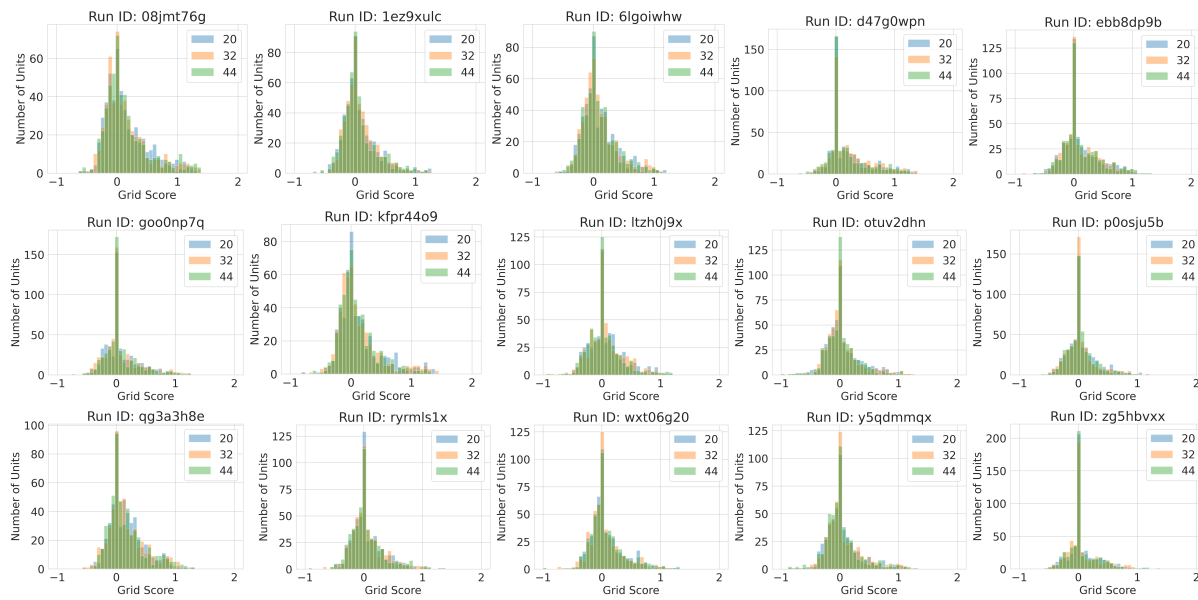


Figure 8: Grid score distributions do not differ as a function of number of bins: 400 (20 x 20; blue), 1024 (32 x 32; orange), 1936 (44 x 44; green).

track used 80 bins.”) but their code similarly used 20 x 20 bins. The code suggests that [36] used the grid scorer of [45], who in turn used the grid scorer of [2].

Consequently, we tested what effect the number of bins has on the distribution of grid scores. We found that the number of bins appears to have little to no effect (Fig. 8), so we chose to use 44 x 44 bins since this yields bins of size 5 cm x 5 cm.

Place cell autocorrelation

In this section, we will derive the form of the place cell correlation function, centered at $\vec{\mu}$. When the spatial tuning curve is a difference-of-Gaussians, the correlation function is *also* a Difference-of-Gaussians, albeit with different parameters. Consider the spatial tuning curve:

$$P(\mathbf{x}; \vec{\mu}) = \alpha_E \exp\left(-\frac{(\mathbf{x} - \vec{\mu})^2}{2\sigma_E^2}\right) - \alpha_I \exp\left(-\frac{(\mathbf{x} - \vec{\mu})^2}{2\sigma_I^2}\right)$$

So the correlation between place cells centered at $\vec{\mu}_1$ and $\vec{\mu}_2$ with $\mathbf{r} = \vec{\mu}_1 - \vec{\mu}_2$ is given by

$$C_{\mu_1, \mu_2} = \int P(\mathbf{x}; \mu_1) P(\mathbf{x}; \mu_2) d\mathbf{x} = \int P(\mathbf{x}; 0) P(\mathbf{x}; \mathbf{r}) d\mathbf{x}$$

Simplifying the above expression, using the identity: $\int dx \mathcal{N}(x; \mu_f, \sigma_f) \mathcal{N}(r-x; \mu_g, \sigma_g) = \mathcal{N}(r; \mu_f + \mu_g, \sqrt{\sigma_f^2 + \sigma_g^2})$, we get:

$$C_{\mu_1, \mu_2} = \alpha_E^2 \exp\left(\frac{-\mathbf{r}^2}{4\sigma_E^2}\right) + \alpha_I^2 \exp\left(\frac{-\mathbf{r}^2}{4\sigma_I^2}\right) - 2\alpha_E \alpha_I \exp\left(\frac{-\mathbf{r}^2}{2(\sigma_E^2 + \sigma_I^2)}\right)$$

Sweeps

.1 Cartesian (Low Dimensional)

```

method: grid
metric:
  goal: minimize
  name: pos_decoding_err
parameters:
  Ng:
    values:
      - 1024
  Np:
    values:
      - 2
  activation:
    values:
      - relu
      - tanh
      - sigmoid
      - linear
  batch_size:
    values:
      - 200
  bin_side_in_m:
    values:
      - 0.05
  box_height_in_m:
    values:
      - 2.2
  box_width_in_m:
    values:
      - 2.2
  initializer:
    values:
      - glorot_uniform
      - glorot_normal
      - orthogonal
  is_periodic:
    values:
      - false
  learning_rate:
    values:
      - 0.0001
  n_epochs:
    values:
      - 20
  n_grad_steps_per_epoch:
    values:
      - 10000
  n_place_fields_per_cell:
    values:
      - 1
  optimizer:
    values:
      - sgd
      - adam

```

```

581         - rmsprop
582     place_cell_rf:
583         values:
584             - 0
585     place_field_loss:
586         values:
587             - mse
588     place_field_normalization:
589         values:
590             - none
591     place_field_values:
592         values:
593             - cartesian
594     readout_dropout:
595         values:
596             - 0
597             - 0.5
598     recurrent_dropout:
599         values:
600             - 0
601             - 0.5
602     rnn_type:
603         values:
604             - RNN
605             - LSTM
606             - UGRNN
607             - GRU
608     seed:
609         values:
610             - 0
611             - 1
612             - 2
613     sequence_length:
614         values:
615             - 20
616     surround_scale:
617         values:
618             - 1
619     weight_decay:
620         values:
621             - 0
622             - 0.0001

```

623 .2 Cartesian (High Dimensional)

```

624 method: grid
625 metric:
626     goal: minimize
627     name: pos_decoding_err
628 parameters:
629     Ng:
630         values:
631             - 1024
632     Np:
633         values:
634             - 512

```

```

635     activation:
636         values:
637             - relu
638             - tanh
639             - sigmoid
640             - linear
641     batch_size:
642         values:
643             - 200
644     bin_side_in_m:
645         values:
646             - 0.05
647     box_height_in_m:
648         values:
649             - 2.2
650     box_width_in_m:
651         values:
652             - 2.2
653     initializer:
654         values:
655             - glorot_uniform
656             - glorot_normal
657             - orthogonal
658     is_periodic:
659         values:
660             - false
661     learning_rate:
662         values:
663             - 0.0001
664     n_epochs:
665         values:
666             - 20
667     n_grad_steps_per_epoch:
668         values:
669             - 1000
670     n_place_fields_per_cell:
671         values:
672             - 1
673     optimizer:
674         values:
675             - adam
676     place_cell_rf:
677         values:
678             - 0
679     place_field_loss:
680         values:
681             - mse
682     place_field_normalization:
683         values:
684             - none
685     place_field_values:
686         values:
687             - high_dim_cartesian
688     readout_dropout:
689         values:
690             - 0

```



```

691         - 0.5
692     recurrent_dropout:
693         values:
694             - 0
695             - 0.5
696     rnn_type:
697         values:
698             - RNN
699             - LSTM
700             - UGRNN
701             - GRU
702     seed:
703         values:
704             - 0
705             - 1
706             - 2
707     sequence_length:
708         values:
709             - 20
710     surround_scale:
711         values:
712             - 0
713     weight_decay:
714         values:
715             - 0
716             - 0.0001

```

717 .3 Polar (High Dimensional)

```

718 method: grid
719 metric:
720     goal: minimize
721     name: pos_decoding_err
722 parameters:
723     Ng:
724         values:
725             - 1024
726     Np:
727         values:
728             - 512
729     activation:
730         values:
731             - relu
732             - tanh
733             - linear
734             - sigmoid
735     batch_size:
736         values:
737             - 200
738     bin_side_in_m:
739         values:
740             - 0.05
741     box_height_in_m:
742         values:
743             - 2.2
744     box_width_in_m:

```

```

745     values:
746         - 2.2
747     initializer:
748         values:
749             - glorot_uniform
750             - glorot_normal
751             - orthogonal
752     is_periodic:
753         values:
754             - false
755     learning_rate:
756         values:
757             - 0.0001
758     n_epochs:
759         values:
760             - 20
761     n_grad_steps_per_epoch:
762         values:
763             - 1000
764     n_place_fields_per_cell:
765         values:
766             - 1
767     optimizer:
768         values:
769             - adam
770     place_cell_rf:
771         values:
772             - 0
773     place_field_loss:
774         values:
775             - mse
776     place_field_normalization:
777         values:
778             - none
779     place_field_values:
780         values:
781             - high_dim_polar
782     readout_dropout:
783         values:
784             - 0
785             - 0.5
786     recurrent_dropout:
787         values:
788             - 0
789             - 0.5
790     rnn_type:
791         values:
792             - RNN
793             - LSTM
794             - UGRNN
795             - GRU
796     seed:
797         values:
798             - 0
799             - 1
800             - 2

```

```

801     sequence_length:
802         values:
803             - 20
804     surround_scale:
805         values:
806             - 0
807     weight_decay:
808         values:
809             - 0
810             - 0.0001

```

811 .4 Gaussian Place Cells

```

812 method: grid
813 metric:
814     goal: minimize
815     name: pos_decoding_err
816 parameters:
817     Ng:
818         values:
819             - 1024
820     Np:
821         values:
822             - 512
823     activation:
824         values:
825             - linear
826             - relu
827             - tanh
828             - sigmoid
829     batch_size:
830         values:
831             - 200
832     bin_side_in_m:
833         values:
834             - 0.05
835     box_height_in_m:
836         values:
837             - 2.2
838     box_width_in_m:
839         values:
840             - 2.2
841     initializer:
842         values:
843             - glorot_uniform
844     is_periodic:
845         values:
846             - false
847     learning_rate:
848         values:
849             - 0.0001
850     n_epochs:
851         values:
852             - 20
853     n_grad_steps_per_epoch:
854         values:

```

```

855         - 10000
856     n_place_fields_per_cell:
857         values:
858             - 1
859     optimizer:
860         values:
861             - adam
862     place_cell_rf:
863         values:
864             - 0.08
865             - 0.1
866             - 0.12
867             - 0.14
868             - 0.16
869             - 0.2
870             - 0.24
871             - 0.28
872     place_field_loss:
873         values:
874             - crossentropy
875     place_field_normalization:
876         values:
877             - global
878     place_field_values:
879         values:
880             - gaussian
881     readout_dropout:
882         values:
883             - 0
884             - 0.5
885     recurrent_dropout:
886         values:
887             - 0
888             - 0.5
889     rnn_type:
890         values:
891             - RNN
892             - LSTM
893             - UGRNN
894             - GRU
895     seed:
896         values:
897             - 0
898             - 1
899     sequence_length:
900         values:
901             - 20
902     surround_scale:
903         values:
904             - 1
905     weight_decay:
906         values:
907             - 0.0001

```

908 .5 Dfference-of-Gaussians Place Cells

```

909 method: grid
910 metric:
911     goal: minimize
912     name: pos_decoding_err
913 parameters:
914     Ng:
915         values:
916             - 1024
917     Np:
918         values:
919             - 512
920     activation:
921         values:
922             - relu
923     batch_size:
924         values:
925             - 200
926     bin_side_in_m:
927         values:
928             - 0.05
929     box_height_in_m:
930         values:
931             - 2.2
932     box_width_in_m:
933         values:
934             - 2.2
935     initializer:
936         values:
937             - glorot_uniform
938     is_periodic:
939         values:
940             - false
941     learning_rate:
942         values:
943             - 0.0001
944     n_epochs:
945         values:
946             - 20
947     n_grad_steps_per_epoch:
948         values:
949             - 10000
950     n_place_fields_per_cell:
951         values:
952             - 1
953     optimizer:
954         values:
955             - adam
956     place_cell_rf:
957         values:
958             - 0.05
959             - 0.1
960             - 0.15
961             - 0.2
962             - 0.3
963             - 0.4
964             - 0.5

```

```

965 place_field_loss:
966     values:
967         - crossentropy
968 place_field_normalization:
969     values:
970         - global
971 place_field_values:
972     values:
973         - true_difference_of_gaussians
974 readout_dropout:
975     values:
976         - 0
977         - 0.5
978 recurrent_dropout:
979     values:
980         - 0
981         - 0.5
982 rnn_type:
983     values:
984         - RNN
985         - LSTM
986         - UGRNN
987         - GRU
988 seed:
989     values:
990         - 0
991         - 1
992         - 2
993 sequence_length:
994     values:
995         - 20
996 surround_scale:
997     values:
998         - 1.5
999         - 2
1000         - 2.5
1001         - 3
1002         - 4
1003         - 5
1004         - 6
1005 weight_decay:
1006     values:
1007         - 0.0001

```

1008 .6 Difference-of-Softmaxes Place Cells

```

1009 method: grid
1010 metric:
1011     goal: minimize
1012     name: pos_decoding_err
1013 parameters:
1014     Ng:
1015         values:
1016             - 1024
1017     Np:
1018         values:

```



```

1019         - 512
1020     activation:
1021         values:
1022             - relu
1023             - tanh
1024     batch_size:
1025         values:
1026             - 200
1027     bin_side_in_m:
1028         values:
1029             - 0.05
1030     box_height_in_m:
1031         values:
1032             - 2.2
1033     box_width_in_m:
1034         values:
1035             - 2.2
1036     initializer:
1037         values:
1038             - glorot_uniform
1039     is_periodic:
1040         values:
1041             - false
1042     learning_rate:
1043         values:
1044             - 0.0001
1045     n_epochs:
1046         values:
1047             - 20
1048     n_grad_steps_per_epoch:
1049         values:
1050             - 10000
1051     n_place_fields_per_cell:
1052         values:
1053             - 1
1054     optimizer:
1055         values:
1056             - adam
1057     place_cell_rf:
1058         values:
1059             - 0.08
1060             - 0.09
1061             - 0.1
1062             - 0.11
1063             - 0.12
1064             - 0.13
1065             - 0.14
1066             - 0.15
1067             - 0.16
1068             - 0.17
1069             - 0.18
1070             - 0.19
1071             - 0.2
1072             - 0.24
1073             - 0.28
1074             - 0.32

```

```

1075 place_field_loss:
1076     values:
1077         - crossentropy
1078 place_field_normalization:
1079     values:
1080         - global
1081 place_field_values:
1082     values:
1083         - difference_of_gaussians
1084 readout_dropout:
1085     values:
1086         - 0
1087         - 0.5
1088 recurrent_dropout:
1089     values:
1090         - 0
1091         - 0.5
1092 rnn_type:
1093     values:
1094         - RNN
1095         - LSTM
1096         - UGRNN
1097         - GRU
1098 seed:
1099     values:
1100         - 0
1101         - 1
1102         - 2
1103 sequence_length:
1104     values:
1105         - 20
1106 surround_scale:
1107     values:
1108         - 1.5
1109         - 2
1110         - 2.5
1111         - 3
1112         - 4
1113 weight_decay:
1114     values:
1115         - 0.0001

```

1116 .7 Softmax-of-Differences Place Cells

```

1117 method: grid
1118 metric:
1119     goal: minimize
1120     name: pos_decoding_err
1121 parameters:
1122     Ng:
1123         values:
1124             - 1024
1125     Np:
1126         values:
1127             - 512
1128     activation:

```

```

1129     values:
1130         - relu
1131     batch_size:
1132         values:
1133             - 200
1134     bin_side_in_m:
1135         values:
1136             - 0.05
1137     box_height_in_m:
1138         values:
1139             - 2.2
1140     box_width_in_m:
1141         values:
1142             - 2.2
1143     initializer:
1144         values:
1145             - glorot_uniform
1146     is_periodic:
1147         values:
1148             - false
1149     learning_rate:
1150         values:
1151             - 0.0001
1152     n_epochs:
1153         values:
1154             - 20
1155     n_grad_steps_per_epoch:
1156         values:
1157             - 10000
1158     n_place_fields_per_cell:
1159         values:
1160             - 1
1161     optimizer:
1162         values:
1163             - adam
1164     place_cell_rf:
1165         values:
1166             - 0.09
1167             - 0.12
1168             - 0.15
1169             - 0.18
1170             - 0.21
1171             - 0.24
1172     place_field_loss:
1173         values:
1174             - crossentropy
1175     place_field_normalization:
1176         values:
1177             - global
1178     place_field_values:
1179         values:
1180             - softmax_of_differences
1181     readout_dropout:
1182         values:
1183             - 0
1184     recurrent_dropout:

```

```

1185     values :
1186         - 0
1187     rnn_type :
1188         values :
1189             - RNN
1190             - LSTM
1191             - UGRNN
1192             - GRU
1193     seed :
1194         values :
1195             - 0
1196             - 1
1197             - 2
1198     sequence_length :
1199         values :
1200             - 20
1201     surround_scale :
1202         values :
1203             - 1.5
1204             - 2
1205             - 2.5
1206             - 3
1207     weight_decay :
1208         values :
1209             - 0.0001

```

1210 .8 Multiple Scales and Multiple Fields Difference-of-Softmaxes Place Cells

```

1211 method: grid
1212 metric :
1213     goal: minimize
1214     name: pos_decoding_err
1215 parameters :
1216     Ng:
1217         values :
1218             - 1024
1219     Np:
1220         values :
1221             - 512
1222     activation :
1223         values :
1224             - relu
1225     batch_size :
1226         values :
1227             - 200
1228     bin_side_in_m :
1229         values :
1230             - 0.05
1231     box_height_in_m :
1232         values :
1233             - 2.2
1234     box_width_in_m :
1235         values :
1236             - 2.2
1237     initializer :
1238         values :

```

```

1239         - glorot_uniform
1240     is_periodic:
1241         values:
1242             - false
1243     learning_rate:
1244         values:
1245             - 0.0001
1246     n_epochs:
1247         values:
1248             - 20
1249     n_grad_steps_per_epoch:
1250         values:
1251             - 10000
1252     n_place_fields_per_cell:
1253         values:
1254             - Poisson( 2.0 )
1255             - Poisson( 3.0 )
1256     optimizer:
1257         values:
1258             - adam
1259     place_cell_rf:
1260         values:
1261             - 0.12
1262             - Uniform( 0.06 , 0.24 )
1263             - Uniform( 0.06 , 1.0 )
1264     place_field_loss:
1265         values:
1266             - crossentropy
1267     place_field_normalization:
1268         values:
1269             - global
1270     place_field_values:
1271         values:
1272             - difference_of_gaussians
1273     readout_dropout:
1274         values:
1275             - 0
1276     recurrent_dropout:
1277         values:
1278             - 0
1279     rnn_type:
1280         values:
1281             - RNN
1282             - UGRNN
1283     seed:
1284         values:
1285             - 0
1286             - 1
1287             - 2
1288     sequence_length:
1289         values:
1290             - 20
1291     surround_scale:
1292         values:
1293             - 2
1294             - Uniform( 1.50 , 2.50 )

```

```

1295         - Uniform( 1.25 , 4.50 )
1296     weight_decay:
1297         values:
1298         - 0.0001

```

1299 .9 Nayebi et al. 2021 [36] Replication

```

1300 method: grid
1301 metric:
1302     goal: minimize
1303     name: pos_decoding_err
1304 parameters:
1305     Ng:
1306         values:
1307         - 4096
1308     Np:
1309         values:
1310         - 512
1311     activation:
1312         values:
1313         - relu
1314     batch_size:
1315         values:
1316         - 200
1317     bin_side_in_m:
1318         values:
1319         - 0.05
1320     box_height_in_m:
1321         values:
1322         - 2.2
1323     box_width_in_m:
1324         values:
1325         - 2.2
1326     initializer:
1327         values:
1328         - glorot_uniform
1329     is_periodic:
1330         values:
1331         - false
1332     learning_rate:
1333         values:
1334         - 0.0001
1335     n_epochs:
1336         values:
1337         - 20
1338     n_grad_steps_per_epoch:
1339         values:
1340         - 10000
1341     n_place_fields_per_cell:
1342         values:
1343         - 1
1344     optimizer:
1345         values:
1346         - adam
1347     place_cell_rf:
1348         values:

```

```

1349         - 0.12
1350 place_field_loss:
1351     values:
1352         - crossentropy
1353 place_field_normalization:
1354     values:
1355         - global
1356 place_field_values:
1357     values:
1358         - difference_of_gaussians
1359 readout_dropout:
1360     values:
1361         - 0
1362 recurrent_dropout:
1363     values:
1364         - 0
1365 rnn_type:
1366     values:
1367         - RNN
1368         - LSTM
1369         - UGRNN
1370         - GRU
1371         - SRNN
1372 seed:
1373     values:
1374         - 0
1375         - 1
1376         - 2
1377         - 3
1378         - 4
1379 sequence_length:
1380     values:
1381         - 20
1382 surround_scale:
1383     values:
1384         - 2
1385 weight_decay:
1386     values:
1387         - 0.0001

```