

Representational drift as a result of implicit regularization

Aviv Ratzon^{1,2*}, Dori Derdikman¹, and Omri Barak^{1,2}

¹*Rappaport Faculty of Medicine, Technion - Israel Institute of Technology, Haifa 31096, Israel*

²*Network Biology Research Laboratory, Technion - Israel Institute of Technology, Haifa 32000, Israel,*

*Aviv.Ratzon@campus.technion.ac.il

Recent studies show that, even in constant environments, the tuning of single neurons changes over time in a variety of brain regions. This representational drift has been suggested to be a consequence of continuous learning under noise, but its properties are still not fully understood. To uncover the underlying mechanism, we trained an artificial network to perform a predictive coding task. After the loss converged, the activity slowly became sparser. We verified the generality of this phenomenon across modeling choices. This sparseness is a manifestation of drift in the solution space to a flatter area. It is consistent with recent experimental results demonstrating that CA1 spatial code becomes sparser after familiarity. We conclude that learning is divided into three overlapping phases: Fast familiarity with the environment, slow implicit regularization, and a steady state of null drift. These findings open the possibility of inferring learning algorithms from observations of drift statistics.

1 What do we mean when we say that the brain represents the external world? One interpretation is the
2 existence of neurons whose activity is tuned to world variables. Such neurons have been observed in
3 many contexts: place cells [1, 2] – which are tuned to position in a specific context, visual cells [3] –
4 which are tuned to specific visual cues, neurons that are tuned to the execution of actions [4] and more.
5 This tight link between the external world and neural activity might suggest that, in the absence of

6 environmental or behavioral changes, neural activity is constant. In contrast, recent studies show that,
7 even in constant environments, the tuning of single neurons to outside world variables gradually changes
8 over time in a variety of brain regions, even long after good representations of the stimuli were achieved.
9 This phenomenon has been termed *representational drift*, and has changed the way we think about the
10 stability of memory and perception, but its driving forces and properties are still unknown [5, 6, 7, 8, 9]
11 (see [10, 11] for an alternative account).

12 There are at least two immediate theoretical questions arising from the observation of drift – why
13 does it happen, and whether and how behavior is resistant to it [12, 13]? One mechanistic explanation
14 is that the underlying anatomical substrates are themselves undergoing constant change, such that drift
15 is a direct manifestation of this structural morphing [14]. A normative interpretation posits that drift
16 is a solution to a computational demand, such as temporal encoding [15], ‘drop-out’ regularization [16],
17 exploration of the solution space [17], or re-encoding during continual learning [12]. Several studies also
18 address the resistance question, providing possible explanations on how behavior can be robust to such
19 phenomena [18, 19, 20, 21].

20 Here, we focus on the mechanistic question, and leverage analyses of drift statistics for this purpose.
21 Specifically, recent studies showed that representational drift in the CA1 is driven by active experience
22 [22]. Namely, rate maps decorrelate more when mice are active for a longer time in a given context.
23 This implies that drift is not just a passive process, but rather an active learning one. As drift seems to
24 occur after an adequate representation has formed, it seems fitting to model it as a form of a continuous
25 learning process.

26 This approach has been recently explored by [23, 24]. They considered continuous learning in noisy,
27 overparameterized neural networks. Because the system is overparameterized, a manifold of zero-loss
28 solutions exists. [23] showed that for feedforward neural networks (FNNs) trained using Hebbian learning
29 with added parameter noise, neurons change their tuning over time. This was due to an *undirected*
30 random walk within the manifold of solutions. The coordinated drift of neighboring place fields was
31 used as evidence to support this view. The phenomenon of undirected motion within the space of
32 solutions seems plausible, as all members of this space achieve equally good performance (Fig 1A left).
33 However, there may be other properties of the solutions (Fig 1B) that vary along this manifold, which
34 could potentially bias drift in a certain direction (Fig 1A right). It is likely that the drift observed in

35 experiments is a combination of both an undirected and directed movement. We will now introduce
36 theoretical results from machine learning that support the possibility of directed drift.

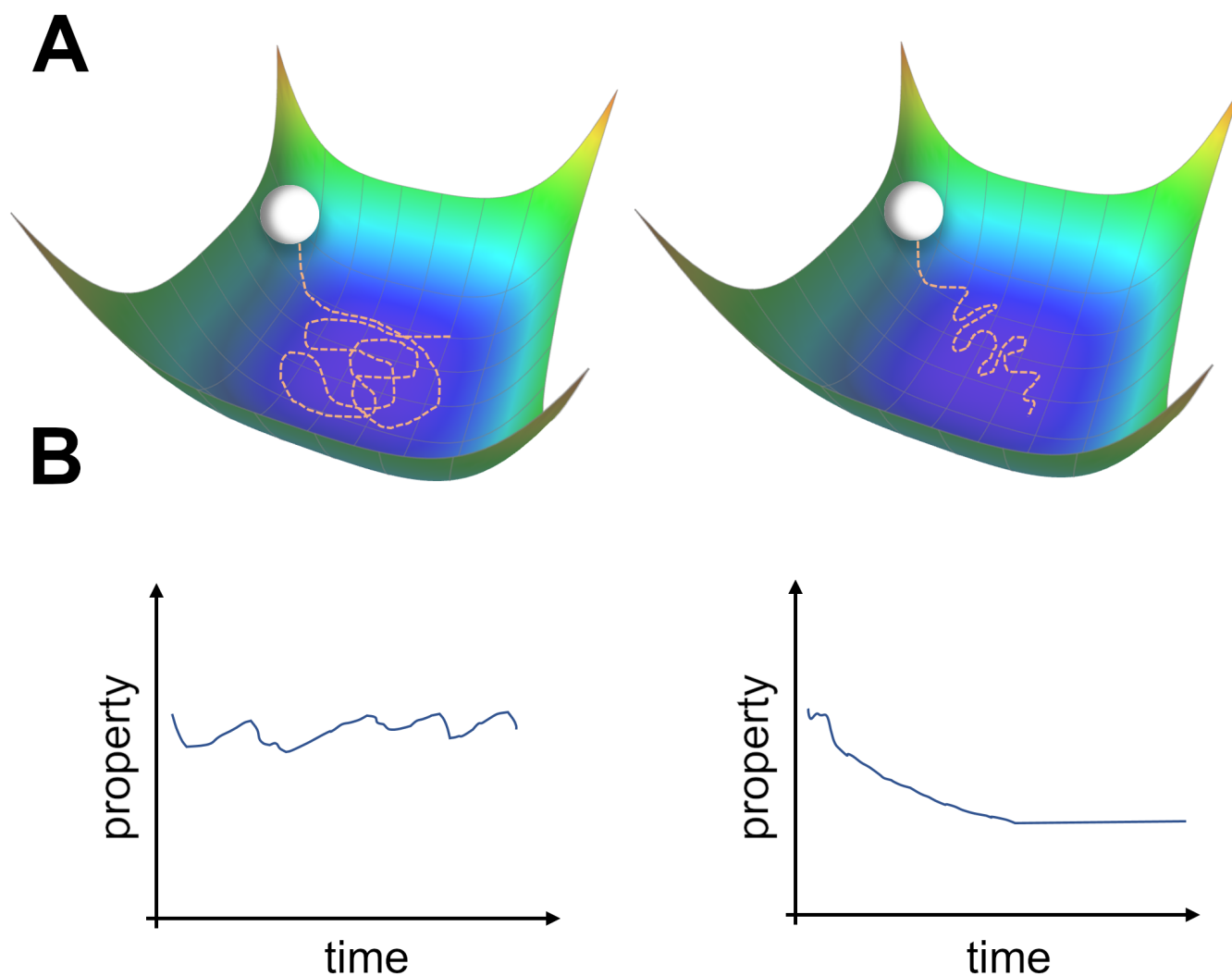


Figure 1: Two types of possible movements within the solution space. (A) Two options of how drift may look in the solution space. Random walk within the space of equally good solutions that is either undirected (left) or directed (right). (B) The qualitative consequence of the two movement types. For an undirected random walk, all properties of the solution will remain roughly constant (left). For the directed movement there should be a given property that is gradually increasing or decreasing (right).

37 Recent work provided a tractable analytical framework for the learning dynamics of Stochastic
38 Gradient Descent (SGD) with added noise and an overparameterized regime [25, 26]. These studies

39 showed that, after the network has converged to the zero-loss manifold, a second-order effect biases
40 the random walk along a specific direction within this manifold. This direction reduces an implicit
41 regularizer, determined by the type of noise the network is exposed to. The regularizer is related to the
42 Hessian of the loss – a measure of the flatness of the loss landscape in the vicinity of the solutions. Since
43 this directed movement is a second-order effect, its timescale is orders of magnitude larger than that of
44 the initial convergence.

45 Consider a biological neural network performing a task. The ML implicit regularization mentioned
46 above requires three components: an overparameterized regime, noise, and SGD. Both biological and
47 artificial networks possess a large number of synapses, or parameters, and hence can reasonably be
48 expected to be overparameterized. Noise can emerge from the external environment or from internal
49 biological elements. It is not reasonable to assume that a precise form of gradient descent is implemented
50 in the brain [27], thereby casting doubt on the third element. Nevertheless, biologically plausible rules
51 could be considered as noisy versions of gradient descent, as long as there is a coherent improvement
52 in performance [28, 29]. Motivated by this analogy, we explore representational drift in models and
53 experimental data.

54 Because drift is commonly observed in spatially-selective cells, we base our analysis on a model
55 which has been shown to contain such cells [30]. Specifically, we trained artificial neural networks on a
56 predictive coding task in the presence of noise. In this task, an agent moves along a linear track while
57 receiving visual input from the walls, such that the goal is to predict the subsequent input. We observed
58 that neurons became tuned to the latent variable, which is position, in accordance with previous results
59 [30]. We continued training and found that in addition to the gradual change of tuning curves, similar
60 to [23], we witnessed that the number of active neurons decreased slowly while their tuning specificity
61 increased. These results align with recent experimental observations [22]. Finally, we demonstrated the
62 connection between this sparsification effect and changes to the Hessian, in accordance with ML theory.

63 Results

64 Spontaneous sparsification in a predictive coding network

65 To model representational drift in the CA1 area, we chose a simple model that could give rise to spatially-
66 tuned cells [30]. In this model, an agent traverses a corridor while slightly modulating its angle with
67 respect to the main axis (Fig 2A). The walls are textured by a fixed smooth noisy signal, and the agent
68 receives this as input according to its current field of view. The model itself is a single hidden layer
69 feedforward network, with the velocity and visual field as inputs. The desired output is the predicted
70 visual input in the next time step. The model equations are given by:

$$\hat{\mathbf{y}}_t = \sigma(\mathbf{x}_t \mathbf{m}^T + \mathbf{b}) \mathbf{n}^T, \quad (1)$$

71 where \mathbf{m} and \mathbf{n} are the input and output matrices respectively, \mathbf{b} is the bias vector, and σ is the ReLU
72 activation function. The task is for the network's output, \mathbf{y} , to match the visual input, \mathbf{x} of the following
73 time step, resulting in the following loss function:

$$f(\mathbf{m}, \mathbf{n}, \mathbf{b}) = \mathbb{E}_t(\hat{\mathbf{y}}_t - \mathbf{x}_{t+1})^2. \quad (2)$$

74 We train the network using Gradient Descent (GD), while adding update noise to the learning
75 dynamics:

$$\theta_{\tau+1} = \theta_{\tau} - \eta \frac{\partial f(\theta_{\tau})}{\partial \theta_{\tau}} + \xi_{\tau}^{update}, \quad (3)$$

76 where $\theta = (\mathbf{m}, \mathbf{n}, \mathbf{b})$ is the vectorized parameters-vector, τ is the current training step and ξ_{τ}^{update} is
77 Gaussian noise. We let the network converge to a good solution, demonstrated by a loss plateau, and
78 continue training for an additional period. Note that this additional period can be orders of magnitude
79 longer than the initial training period. The network quickly converged to a low loss and stayed at the
80 same loss during the additional training period (Fig 2B). Surprisingly, when looking at the activity within
81 the hidden layer, we noticed that it slowly became sparse. This sparsification did not hurt performance,

82 because individual units became more informative, as quantified by the average mutual information
83 between unit activity and the position of the agent (Fig 2C). When looking at the rate maps of neurons,
84 i.e. their tuning to position, one can observe an image similar to representational drift observed in
85 experiments [5] – namely that neurons changed their tuning over time (Fig 2D). Additionally, their
86 tuning specificity increased in accordance with the information increase. By observing the correlation
87 matrix of the rate maps over time, it is apparent that there was a gradual change that slowed down
88 (Fig 2E). To summarize, we observed a spontaneous sparsification over a timescale much longer than the
89 initial convergence, without introducing any explicit regularization. This is comparable to experimental
90 data from [22], where indeed drift was characterized by a decrease in the fraction of active place cells, and
91 an increase in cells’ information while the decoding error for the position of the mouse stayed relatively
92 constant (Fig 2F). Another recent study further demonstrated an increase in information over days [31].

93 **Generality of the phenomenon**

94 To explore the sensitivity of our results to specific modeling choices, we systematically varied many of
95 them (Fig 3A). Specifically, we replaced the task with either a simplified predictive coding, random
96 mappings or smoothed random mappings. Noise was introduced to the outputs (label noise), instead
97 of the update noise. We simulated different activation functions. Perhaps most important, we varied
98 the learning rules, as SGD is not a biologically plausible one. We used both Adam [32] and RMSprop
99 [33], from the ML literature. We also used Stochastic Error-Descent (SED) [34], which does not require
100 gradient calculation and is more biologically plausible (5). All cases demonstrated an initial, fast, phase
101 of convergence to low loss, followed by a much slower phase of directed random motion within the
102 low-loss space.

103 The results of the simulations supported our main conclusion, though several qualitative phe-
104 nomenons could be observed. First of all, sparsification dynamics were not sensitive to most of the
105 parameters. The main qualitative difference observed was that the timescales could vary by orders of
106 magnitude as a function of the noise scale (Fig 3B bottom). Note that we calculate the timescale of
107 sparsification by fitting an exponential curve to the fraction of active units over time, and take the time
108 constant of the fitted exponential. Additionally, apart from simulations that did not converge due to
109 too big timescales, the final sparsity was the same for all networks of the same size (Fig 3B top), in

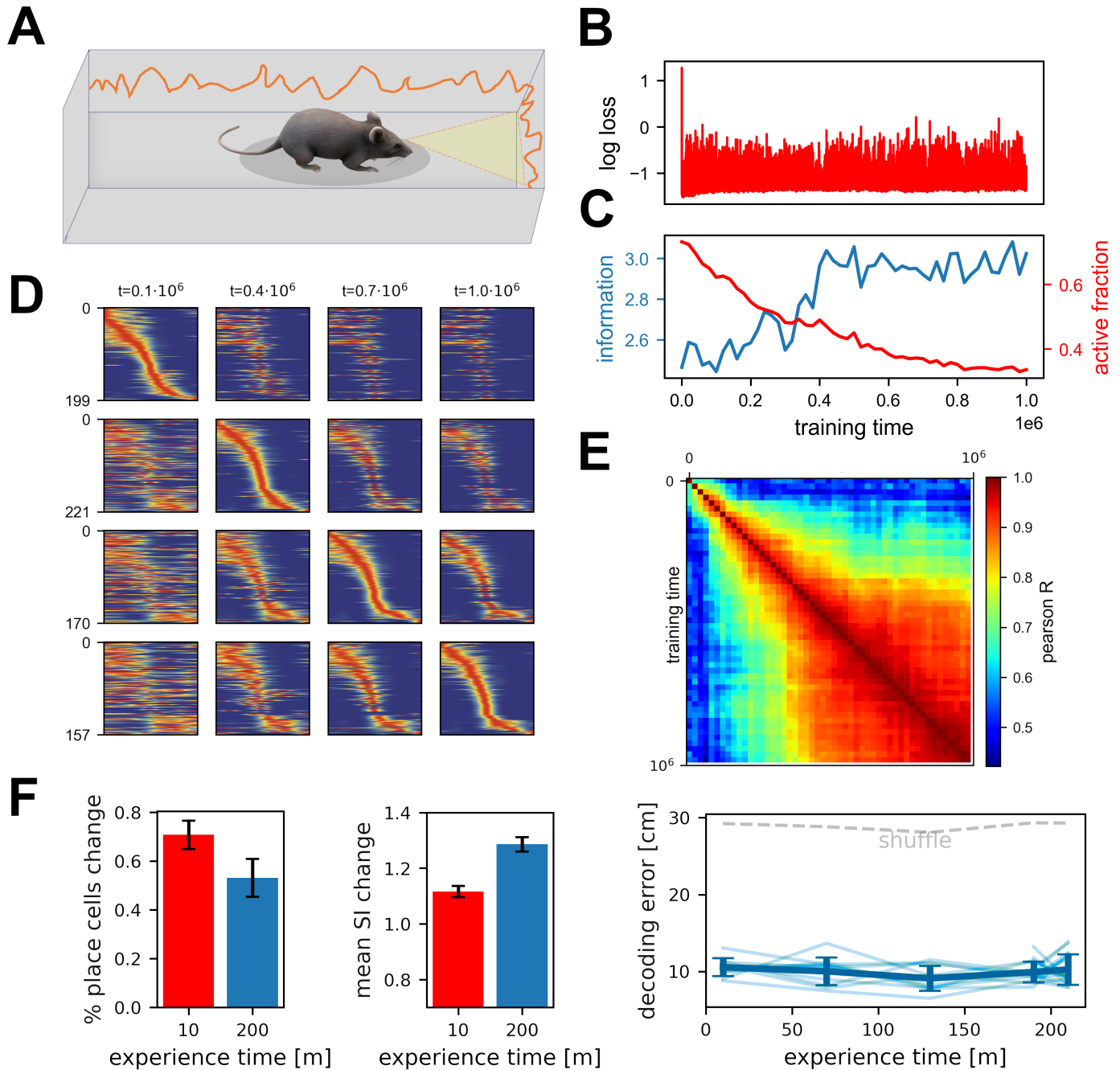


Figure 2: Noisy learning leads to spontaneous sparsification. (A) Illustration of an agent in a corridor receiving high-dimensional visual input from the walls. (B) Log loss as a function of training steps, log loss of 0 corresponds to a mean estimator. The loss rapidly decreases, and then remains roughly constant. (C) Information (blue) and fraction of units with non-zero activation for at least one input (red) as a function of training steps. (D) Rate maps sampled at four different time points. Maps in each row are sorted according to a different time point. Sorting is done based on the peak tuning value to the latent variable. (E) Correlation of rate maps between different time points along training. Only active units are used. (F) Figures reproduced from [22] where mice spent different amount of time in two environments. Fraction of place cells in the beginning relative to the end of the experiment (left), average Spatial Information (SI) per cell in the beginning relative to the end of the experiment (middle) and the decoding error for the position of the mouse (right).

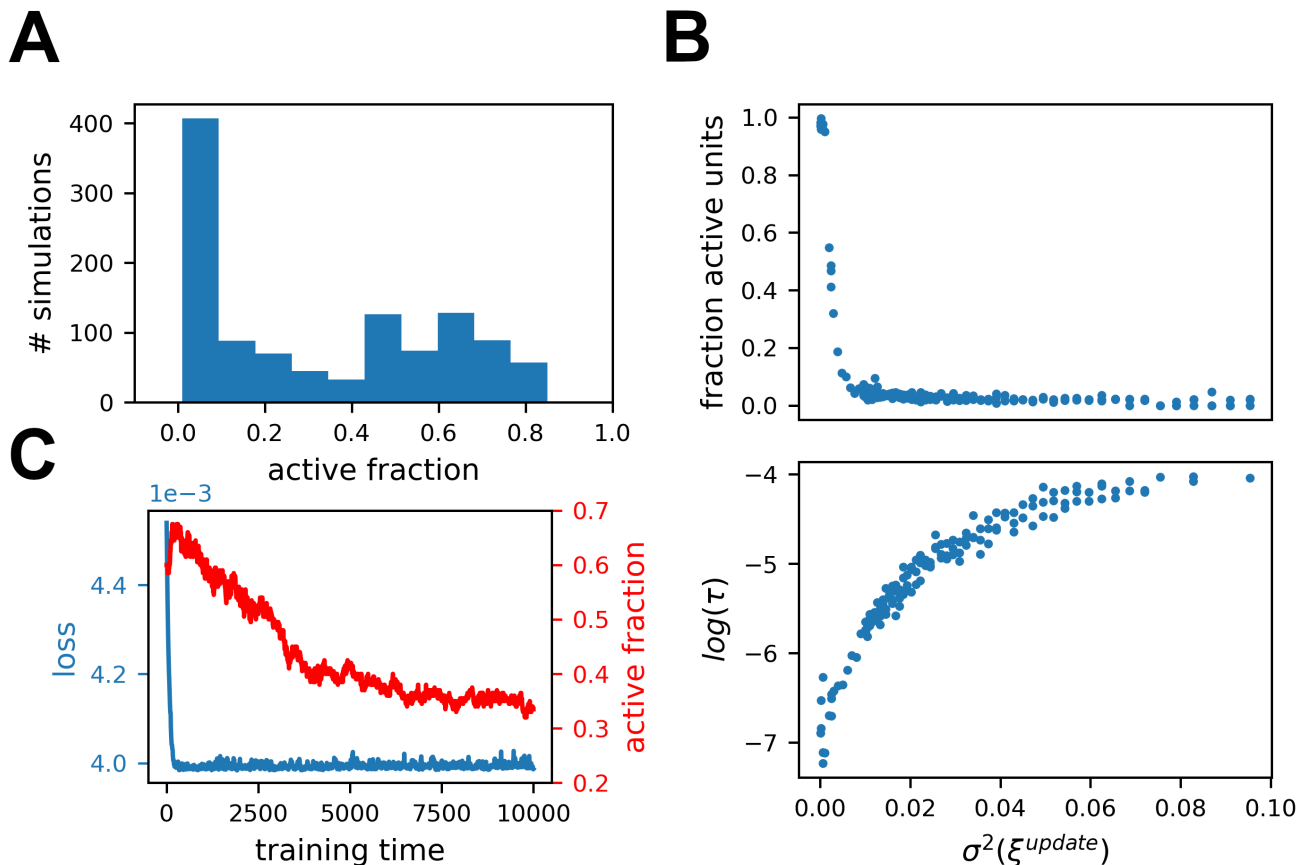


Figure 3: Generality of the results. Summary of 1117 simulations with various parameters (see Table 1). (A) Histogram of fraction of active units after 10^7 training steps for each simulation. (B) Subset of 178 simulations with the same parameters and varying noise variance, each point represents a single simulation. Fraction of active units as a function of the variance of the noise (top), the log of sparsification time scale as a function of the variance of the noise (bottom). (C) Learning a similarity matching task with Hebbian and anti-Hebbian learning with published code from [23]. Performance of the network (blue) and fraction of active units (red) as a function of training steps. Note that the loss axis does not start at zero, and the dynamic range is small.

110 accordance with results from [23]. In a sense, once noise is introduced the network is driven to maximal
 111 sparsification. For Adam, RMSprop and SED sparsification ensued in the absence of any added noise.
 112 For SED the explanation is straightforward, as the parameter updates are driven by noise. For Adam
 113 and RMSprop, we suggest that in the vicinity of the zero-loss manifold, the second moment acts as noise.
 114 For label noise, the dynamics were qualitatively different, the fraction of active units did not reduce,
 115 but the activity of the units did sparsify. In some cases, the networks quickly collapsed to a sparse

116 solution, most likely as a result of the learning rate being too high, in relation to the input statistics
117 [35]. Importantly, for GD without noise, there was no change after the initial convergence.

118 As a further test of the generality of this phenomenon, we consider the recent simulation from [23].
119 The learning rule used in this work was very different from the ones we applied. We, therefore, simulated
120 that network using the published code. We found the same type of dynamics as shown above, namely that
121 the network initially converged to a good solution followed by a longer period of sparsification (Fig 3C).
122 Note that in The original publication [23] the focus was on the stage following this sparsification, in
123 which the network indeed maintained a constant fraction of active cells.

124 In conclusion, we see that noisy learning leads to three phases under rather general conditions. First,
125 fast learning of the task and convergence to the manifold of low-loss solutions. The second phase is
126 directed movement on this manifold driven by a second-order effect of implicit regularization. The third
127 phase is an undirected random walk within the sub-manifold of low loss and maximum regularization.

128 **Mechanism of sparsification**

129 What are the mechanisms that give rise to sparsification? As illustrated in Fig. 1, different solutions
130 in the zero-loss manifold might vary in some of their properties. The specific property suggested from
131 theory [25] is the flatness of the loss landscape in the vicinity of the solution. This can be demonstrated
132 with a simple example. Consider a two-dimensional loss function. The function is shaped like a valley
133 with a continuous one-dimensional zero-loss manifold at its bottom (Fig 4A). Crucially, the loss on the
134 entire manifold is exactly zero, while the vicinity of the manifold becomes systematically flatter in one
135 direction. We simulated gradient descent with added noise on this function from a random starting
136 point (red dot). The trajectory quickly converged to the zero-loss manifold, and began a random walk
137 on it. This walk was clearly biased towards the flatter area of the manifold, as can be seen by the spread
138 of the trajectory. This bias could be comprehended by noting that the gradient was orthogonal to the
139 contour lines of the loss, and therefore had a component directed towards the flat region.

140 In higher dimensions, flatness is captured by the eigenvalues of the Hessian of the loss. Because these
141 eigenvalues are a collection of numbers, different scenarios could lead to minimizing different aspects
142 of this collection. Specifically, according to [25], update noise should regularize the sum of the log
143 of the non-zero eigenvalues while label noise should do the same for the sum of eigenvalues. In our

144 predictive coding example, where update noise was added, each inactivated unit translates into a set
145 of zero-rows in the Hessian, and thus also into a set of zero-eigenvalues (Fig 4B). The slope of the
146 regularizer approaches infinity as the eigenvalue approaches zero, and thus small eigenvalues are driven
147 to zero much faster than large eigenvalues (Fig 4C). So in this case, update noise leads to an increase
148 in the number of zero eigenvalues, which are manifested as a sparse solution. Another, perhaps more
149 intuitive, way to understand these results is that units below the activation threshold are insensitive to
150 noise perturbations. In other scenarios, in which we simulated with label noise, we indeed observed a
151 gradual decrease in the sum of eigenvalues (Fig 4D).

152 Discussion

153 We showed that representational drift could arise from ongoing learning in the presence of noise, af-
154 ter a network has already reached good performance. We suggest that learning is divided into three
155 overlapping phases: a fast initial phase, where good performance is achieved, a second slower phase in
156 which *directed* drift along the low-loss manifold leads to an implicit regularization and finally, a third
157 *undirected* phase ensues once the regularizer is minimized. In our results, the directed component was
158 associated with sparsification of the neural code, a phenomenon we also observed in experimental data.

159 Interpreting drift as a learning process has recently been suggested by [23, 24]. Both studies focused
160 on the final phase in which the statistics of the representations were constant. Experimentally, [7]
161 reported a decrease in activity at the beginning of the experiment, which they suggested was correlated
162 with some behavioral change, but we believe it could also be a result of the directed drift phase. [36]
163 also reported a slow directed change in representation long after familiarity with the stimuli. There
164 is another consequence of the timescale separation. Unlike in the setting of drift experiments, natural
165 environments are never truly constant. Thus, it is possible that the second phase of learning never stops
166 because the task is slowly changing. This would imply that the second, directed, phase may be the
167 natural regime in which neural networks reside.

168 Here, we reported directed drift in the space of solutions of neural networks. This drift could be
169 observed by examining changes to the representation of external world variables, and hence is related to
170 the phenomenon of representational drift. Note, however, that representations are not a full description

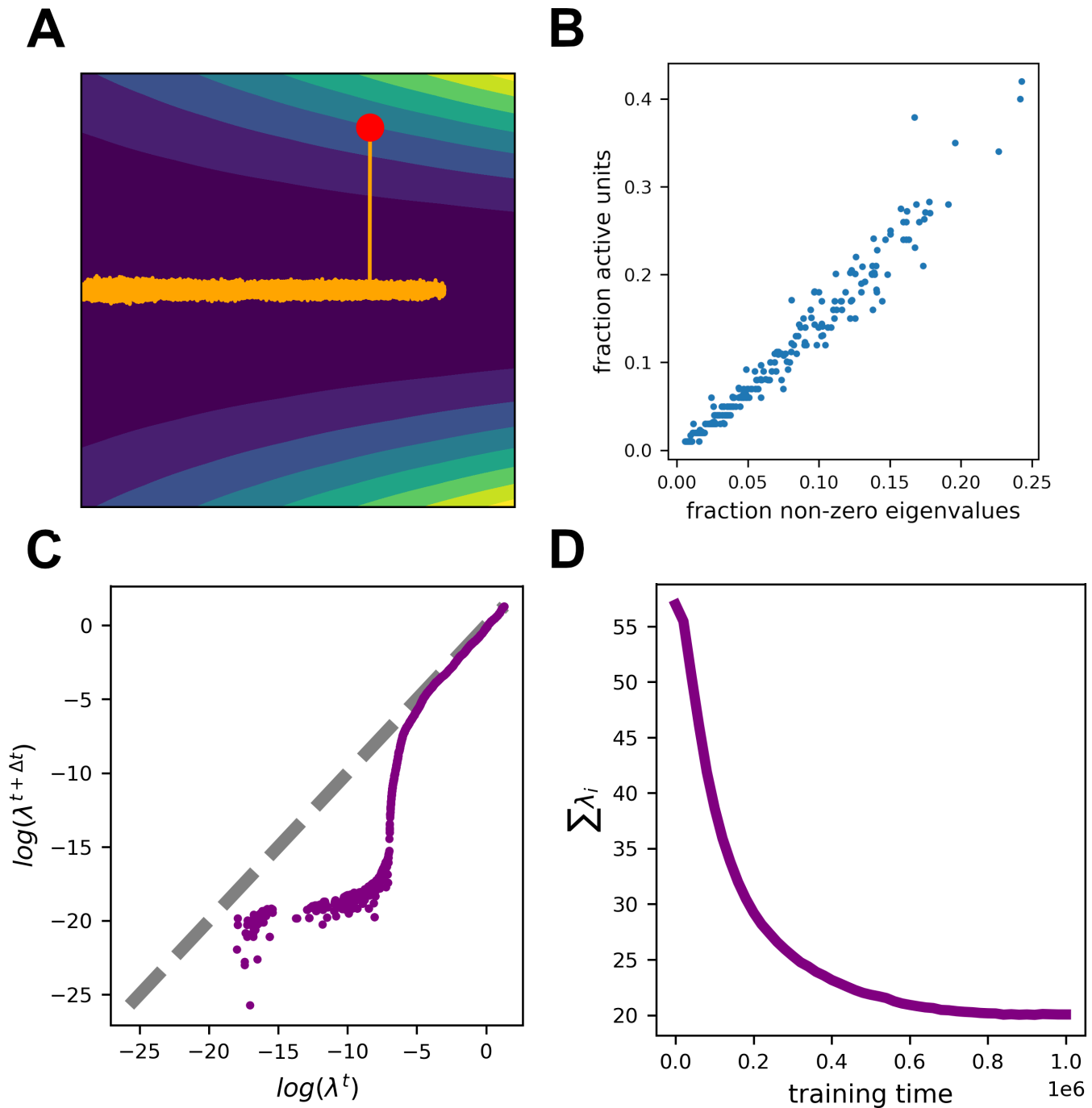


Figure 4: Noisy learning leads to a flat landscape. (A) Gradient Descent dynamics over a two-dimensional loss function with a one-dimensional zero-loss manifold. Note that the loss is identically zero along the horizontal axis, but the left area is flatter. The orange trajectory begins at the red dot. Note the asymmetric extension into the left area. (B) Fraction of active units as a function of the number of non-zero eigenvalues. (C) Log of non-zero eigenvalues at two consecutive time points. Note that eigenvalues do not correspond to one another when calculated at two different time points, and this plot demonstrates the change in their distribution rather than changes in eigenvalues corresponding to specific directions. The distribution of larger eigenvalues hardly changes, while the distribution of smaller eigenvalues is pushed to smaller values. (D) Sum of the Hessian's eigenvalues as a function of time for learning with label noise.

171 of a network’s behavior [37]. The statistics of representational changes can be used as a window into
172 changes of network dynamics and function.

173 The phenomenon of directed drift is very robust to various modeling choices, and also consistent
174 with recent theoretical results [25, 26] The details of the direction of the drift, however, are dependent
175 on specific choices. Specifically, which aspects of the Hessian are minimized during the second phase of
176 learning, as well as the timescale of this phase, depend on the specifics of the learning rule and the noise
177 in the system. This suggests an exciting opportunity – inferring the learning rule of a network from the
178 statistics of representational drift.

179 Our explanation of drift invoked the concept of a low-loss manifold – a family of network config-
180 urations that have identical performance on a task. The definition of low-loss, however, depends on
181 the specific task and context analyzed. Challenging a system with new inputs could dissociate two
182 configurations that otherwise appear identical [38]. It will be interesting to explore whether various en-
183 vironmental perturbations could uncover the motion along the low-loss manifold in the CA1 population.
184 For instance, remapping was interpreted as an indicator of the detection of a context switch [39]. One
185 can therefore speculate that the probability for remapping given the same environmental change will
186 systematically vary as the network moves to flatter areas of the loss landscape.

187 Machine learning has been suggested as a model tool for neuroscience research [40, 41, 42]. However,
188 the implicit regularization in ML has not been studied to explain representational drift in neuroscience,
189 and may have been done without awareness of this phenomenon. It’s worth noting that this isn’t a
190 phenomenon specific to neural networks, but rather a general property of overparameterized systems
191 that optimize a cost function. Importing insights from this domain into neuroscience shows the utility
192 of studying general phenomena in systems that learn. For example, another complex learning system in
193 which a similar idea has been proposed is evolution – ”survival of the flattest” suggests that, under a
194 high mutation rate, the fittest replicators are not just the ones with the highest fitness, but also with a
195 flat fitness function which is more robust to mutations [43]. One can hope that more such insights will
196 arise as we open our eyes.

197 **Materials and methods**

198 **Predictive coding task**

199 The agent is moving in an arena of size (L_x, L_y) , with constant velocity in the y direction of V_0 . The
200 agent's heading direction is θ and it changes at every time step by $\Delta\theta \sim G(0, \sigma_\theta^2)$, the agent's visual
201 field has an angle θ_{vis} and is represented as a vector of size L_{vis} . The texture of the walls is generated
202 from a random Gaussian vector of size $L_{walls} = 2(L_x + L_y)L_{vis}$, smoothed with a Gaussian filter with
203 $\sigma^2 = K_{smooth}L_{walls}$. At each time step the agent receives the visual input from the walls, determined by
204 the intersection points of it's visual field with the walls. When the agent reaches a distance of L_yL_{buffer}
205 from the wall, it turns to the opposite direction.

206 **Tuning properties of units**

207 For each unit we calculated a tuning curve. We divided the arena into 100 equal bins and computed the
208 number of time steps in each bin and the mean unit activation. We then obtained the tuning curve by
209 dividing the mean activity for each bin by the occupancy. We treated movement in each direction as a
210 separate location. We calculated the spatial information (SI) of the tuning curves for each unit:

$$SI = \sum_i p_i \frac{r_i}{\bar{r}} \log_2 \frac{r_i}{\bar{r}} \quad (4)$$

211 where i is the index of the bin, p_i is the probability of being in the bin, r_i is the value of the tuning
212 curve in the bin and \bar{r} is the unit's mean activity rate. Active unit was defined as a unit with non-zero
213 activation for at least one input.

214 **Simulations**

215 For the random simulations, we train each network for 10^7 training steps while choosing random learning
216 algorithm and parameters. The ranges and relevant values of parameters are specified in Table 1. For
217 Adam and SED there was no added noise.

Table 1: Parameter ranges for random simulations.

Parameter	Possible values
learning algorithm	{SGD, Adam, SED}
noise type	{update, label}
number of samples	[1,20]
initialization regime	{lazy, rich}
task	{abstract predictive, random, random smoothed}
input dimension	[1,20]
output dimension	[1,20]
noise variance (label/update)	[0.1,1]/[0.01,0.1]
hidden layer size	100

218 **Stochastic Error Descent**

219 The equation for parameter updates under this learning rule is given by:

$$\theta_{\tau+1} = \theta_{\tau} - \eta(f(\theta_{\tau} + \xi_{\tau}) - f(\theta_{\tau}))\xi_{\tau} \quad (5)$$

220 In this learning rule, the parameters are randomly perturbed at each training step by a Gaussian
221 noise denoted by ξ_{τ} and then updated in proportion to the change in loss.

222 **Label noise**

223 Label noise is introduced to the loss function given by the following formula:

$$f(\mathbf{x}_t) = (\hat{\mathbf{y}}_t - \mathbf{x}_{t+1} + \xi_{\tau}^{label})^2, \quad (6)$$

224 where ξ_{τ}^{label} is Gaussian noise.

225 **Gradient descent dynamics around the zero-loss manifold**

226 The function we used for the two-dimensional example was given by:

$$L(x, y) = (xy)^2, \quad (7)$$

227 which has zero loss on the x and y axes. For small enough update noise, GD will converge to the vicinity
228 of this manifold (the axes). We consider a point on the x axis: $(x_0, 0)$, and calculate the direction of the
229 gradient near that point. Because we are interested in motion along the zero-loss manifold, we consider
230 a small perturbation in the orthogonal direction $(x_0, 0 + \Delta y)$ where $x_0 \gg 1$ and $|\Delta y| \ll 1$. Any
231 component of the gradient in the x direction will lead to motion along the manifold. The update step
232 at this point is given by:

$$-\nabla L(x_0, 0 + \Delta y) = -2x_0 \begin{pmatrix} (\Delta y)^2 \\ x_0 \Delta y \end{pmatrix}. \quad (8)$$

233 One can observe that the step has a large component in the y direction, quickly returning to the
234 manifold. There is also a smaller component in the x direction, reducing the value of x . Reducing x
235 also reduces the Hessian's eigenvalues:

$$H_L(x_0, 0) = 2 \begin{pmatrix} 0 & 0 \\ 0 & x_0^2 \end{pmatrix} \quad (9)$$

$$\lambda_{1,2} = \{0, x_0^2\}, v_{1,2} = \left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right\}. \quad (10)$$

236 Thus, it becomes clear that the trajectory will have a bias that reduces the curvature in the y
237 direction.

238 For general loss functions and various noise models, rigorous proofs can be found in [25], and a
239 different approach can be found in [26]. Here, we will briefly outline the intuition for the general case.
240 Consider again the update rule for GD:

$$\theta \leftarrow \theta - \eta \nabla L(\theta). \quad (11)$$

241 In order to understand the dynamics close to the zero-loss manifold, we consider a point θ , for which
242 $L(\theta) = 0$ expand the loss around it:

$$L(\theta + \delta\theta) = L(\theta) + \nabla^T L(\theta)\delta\theta + \frac{1}{2}\delta\theta^T H\delta\theta. \quad (12)$$

243 We can then take the gradient of this expansion with respect to θ :

$$\nabla_{\theta}L(\theta + \delta\theta) = \nabla_{\theta}L(\theta) + \nabla_{\theta}\nabla_{\theta}^T L(\theta)\delta\theta + \nabla_{\theta}\left(\frac{1}{2}\delta\theta^T H\delta\theta\right) \quad (13)$$

$$= 0 + H\delta\theta + \nabla_{\theta}\left(\frac{1}{2}\delta\theta^T H\delta\theta\right). \quad (14)$$

244 The first term is zero, because the gradient is zero on the manifold. The second term is the largest
245 one, as it linear in $\delta\theta$. Note that the Hessian matrix has zero eigenvalues in directions on the zero-loss
246 manifold, and non-zero eigenvalues in other directions. Thus, the second term corresponds to projecting
247 $\delta\theta$ in a direction that is orthogonal to the zero-loss manifold. The third term can be interpreted as
248 the gradient of some auxiliary loss function. Thus, we expect gradient descent to minimize this new
249 loss, which corresponds to a quadratic form with the Hessian. This is the reason for the implicit
250 regularization along the manifold. Note that the auxiliary loss function is defined by $\delta\theta$, and thus
251 different noise statistics will correspond, on average, to different implicit regularizations. In conclusion,
252 the update step will have a large component that moves the parameter vector towards the zero-loss
253 manifold, and a small component that moves the parameter vector on the manifold in a direction that
254 minimizes some measure of the Hessian.

255 Hessian and sparseness

256 In the main text, we show that the implicit regularization of the Hessian leads to sparse representa-
257 tions. Here, we show this relationship for a single-hidden layer feed-forward neural network with ReLU
258 activation and Mean Squared Error loss:

$$f(\mathbf{x}_i) = \sigma(\mathbf{x}_i \mathbf{m}^T + b) \mathbf{n}^T \quad (15)$$

259 The gradient and Hessian at the zero-loss manifold are given by [44]:

$$\nabla_{\theta} f(\mathbf{x}_i) = \begin{pmatrix} \frac{\partial f}{\partial \mathbf{m}} \\ \frac{\partial f}{\partial \mathbf{b}} \\ \frac{\partial f}{\partial \mathbf{n}} \end{pmatrix} = \begin{pmatrix} \mathbf{n} \odot \mathbb{1}(\mathbf{x}_i; \theta) \otimes \mathbf{x}_i \\ \mathbf{n} \odot \mathbb{1}(\mathbf{x}_i; \theta) \\ (\mathbf{x}_i \cdot \mathbf{n}^T + \mathbf{b}) \odot \mathbb{1}(\mathbf{x}_i; \theta) \end{pmatrix} \quad (16)$$

$$\nabla_{\theta}^2 L(\mathbf{x}; \theta) = \sum_i \nabla_{\theta} f(\mathbf{x}_i) \nabla_{\theta} f(\mathbf{x}_i)^T, \quad (17)$$

260 where $\mathbb{1}(\mathbf{x}_i; \theta)$ is an indicator vector denoting whether each unit is active for some input \mathbf{x}_i . Sparseness
 261 means that a unit has become inactive for all inputs. All the partial derivatives of input, output and
 262 bias weights associated with such a unit are zero, and thus the relevant rows of the Hessian are zero as
 263 well. Thus, every inactive unit leads to several zero eigenvalues.

264 Acknowledgments

265 We thank Ron Teichner and Kabir Dabholkar for comments on the manuscript. This research was
 266 supported by the ISRAEL SCIENCE FOUNDATION (grants Nos. 2655/18 and 2183/21 to DD, and
 267 1442/21 to OB), by the German-Israeli Foundation (GIF I-1477-421.13/2018) to DD, by a grant from
 268 the US-Israel Binational Science Foundation (NIMH-BSF CRCNS BSF:2019807, NIMH:R01 MH125544-
 269 01 to DD), by an HFSP research grant (RGP0017/2021) to OB, A Rappaport Institute Collaborative
 270 research grant to DD, by Israel PBC-VATAT and by the Technion Center for Machine Learning and
 271 Intelligent Systems (MLIS) to DD and OB, by the Prince Center for the Aging Brain, and by a University
 272 of Michigan – Israel Partnership for Research and Education Collaborative Research stipend to DK.
 273 (data science)

References

- [1] John O’keefe and Lynn Nadel. The hippocampus as a cognitive map. *Behavioral and Brain Sciences*, 2(4):487–494, 1979.
- [2] John O’Keefe and Jonathan Dostrovsky. The hippocampus as a spatial map: preliminary evidence from unit activity in the freely-moving rat. *Brain research*, 1971.

- [3] David H Hubel and Torsten N Wiesel. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *The Journal of physiology*, 160(1):106, 1962.
- [4] Bruce L McNaughton, SJY Mizumori, CA Barnes, BJ Leonard, M Marquis, and EJ Green. Cortical representation of motion during unrestrained spatial navigation in the rat. *Cerebral Cortex*, 4(1):27–39, 1994.
- [5] Yaniv Ziv, Laurie D Burns, Eric D Cocker, Elizabeth O Hamel, Kunal K Ghosh, Lacey J Kitch, Abbas El Gamal, and Mark J Schnitzer. Long-term dynamics of ca1 hippocampal place codes. *Nature neuroscience*, 16(3):264–266, 2013.
- [6] Laura N. Driscoll, Noah L. Pettit, Matthias Minderer, Selmaan N. Chettih, and Christopher D. Harvey. Dynamic Reorganization of Neuronal Activity Patterns in Parietal Cortex. *Cell*, 170(5):986–999, 8 2017.
- [7] Daniel Deitch, Alon Rubin, and Yaniv Ziv. Representational drift in the mouse visual cortex. *bioRxiv*, page 2020.10.05.327049, 10 2020.
- [8] Carl E Schoonover, Sarah N Ohashi, Richard Axel, and Andrew J P Fink. Representational drift in primary olfactory cortex. *Nature*, 594, 2021.
- [9] Ziv Y. Geva N., Rubin A. Differential effects of time and experience on hippocampal representational drift. *[Poster]*. In: *Cosyne Convention, 2022, March 17-20, Lisbon, Portugal*, 2022.
- [10] William A Liberti, 3rd, Tobias A Schmid, Angelo Forli, Madeleine Snyder, and Michael M Yartsev. Publisher correction: A stable hippocampal code in freely flying bats. *Nature*, 606(7914):E6, June 2022.
- [11] Sadra Sadeh and Claudia Clopath. Contribution of behavioural variability to representational drift. *Elife*, 11:e77907, 2022.
- [12] Michael E Rule, Timothy O’Leary, and Christopher D Harvey. Causes and consequences of representational drift. *Curr. Opin. Neurobiol.*, 58:141–147, October 2019.
- [13] Laura N Driscoll, Lea Duncker, and Christopher D Harvey. Representational drift: Emerging theories for continual learning and experimental future directions. *Current Opinion in Neurobiology*, 76:102609, 2022.
- [14] Noam E Ziv and Naama Brenner. Synaptic tenacity or lack thereof: spontaneous remodeling of synapses. *Trends in neurosciences*, 41(2):89–99, 2018.

- [15] Alon Rubin, Nitzan Geva, Liron Sheintuch, and Yaniv Ziv. Hippocampal ensemble dynamics timestamp events in long-term memory. *elife*, 4:e12247, 2015.
- [16] Kyle Aitken, Marina Garrett, Shawn Olsen, and Stefan Mihalas. The geometry of representational drift in natural and artificial neural networks. *PLOS Computational Biology*, 18(11):e1010716, 2022.
- [17] David Kappel, Stefan Habenschuss, Robert Legenstein, and Wolfgang Maass. Network plasticity as bayesian inference. *PLoS computational biology*, 11(11):e1004485, 2015.
- [18] Uri Rokni, Andrew G Richardson, Emilio Bizzi, and H Sebastian Seung. Motor learning with unstable neural representations. *Neuron*, 54(4):653–666, 2007.
- [19] Lee Susman, Naama Brenner, and Omri Barak. Stable memory with unstable synapses. *Nature communications*, 10(1):4441, 2019.
- [20] Gianluigi Mongillo, Simon Rumpel, and Yonatan Loewenstein. Intrinsic volatility of synaptic connections—a challenge to the synaptic trace theory of memory. *Current opinion in neurobiology*, 46:7–13, 2017.
- [21] Yaroslav Felipe Kalle Kossio, Sven Goedeke, Christian Klos, and Raoul-Martin Memmesheimer. Drifting assemblies for persistent memory: Neuron transitions and unsupervised compensation. *Proceedings of the National Academy of Sciences*, 118(46):e2023832118, 2021.
- [22] Dorgham Khatib, Aviv Ratzon, Mariell Sellevoll, Omri Barak, Genela Morris, and Dori Derdikman. Experience, not time, determines representational drift in the hippocampus. *bioRxiv*, pages 2022–08, 2022.
- [23] Shanshan Qin, Shiva Farashahi, David Lipshutz, Anirvan M Sengupta, Dmitri B Chklovskii, and Cengiz Pehlevan. Coordinated drift of receptive fields in hebbian/anti-hebbian network models during noisy representation learning. *Nature Neuroscience*, pages 1–11, 2023.
- [24] Farhad Pashakhanloo and Alexei Koutrakov. Stochastic gradient descent-induced drift of representation in a two-layer neural network. *arXiv preprint arXiv:2302.02563*, 2023.
- [25] Guy Blanc, Neha Gupta, Gregory Valiant, and Paul Valiant. Implicit regularization for deep neural networks driven by an ornstein-uhlenbeck like process. In *Conference on learning theory*, pages 483–513. PMLR, 2020.
- [26] Zhiyuan Li, Tianhao Wang, and Sanjeev Arora. What happens after sgd reaches zero loss?—a mathematical framework. *arXiv preprint arXiv:2110.06914*, 2021.

- [27] Yoshua Bengio, Dong-Hyun Lee, Jorg Bornschein, Thomas Mesnard, and Zhouhan Lin. Towards biologically plausible deep learning. *arXiv preprint arXiv:1502.04156*, 2015.
- [28] Yuhan Helena Liu, Arna Ghosh, Blake Richards, Eric Shea-Brown, and Guillaume Lajoie. Beyond accuracy: generalization properties of bio-plausible temporal credit assignment rules. *Advances in Neural Information Processing Systems*, 35:23077–23097, 2022.
- [29] Owen Marschall, Kyunghyun Cho, and Cristina Savin. A unified framework of online learning algorithms for training recurrent neural networks. *The Journal of Machine Learning Research*, 21(1):5320–5353, 2020.
- [30] Stefano Recanatesi, Matthew Farrell, Guillaume Lajoie, Sophie Deneve, Mattia Rigotti, and Eric Shea-Brown. Predictive learning as a network mechanism for extracting low-dimensional latent space representations. *Nature Communications*, 12(1), 2021.
- [31] Liron Sheintuch, Alon Rubin, and Yaniv Ziv. Bias-free estimation of information content in temporally sparse neuronal activity. *PLoS computational biology*, 18(2):e1009832, 2022.
- [32] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [33] Geoffrey Hinton, Nitish Srivastava, and Kevin Swersky. Neural networks for machine learning lecture 6a overview of mini-batch gradient descent. *Cited on*, 14(8):2, 2012.
- [34] Gert Cauwenberghs. A fast stochastic error-descent algorithm for supervised learning and optimization. *Advances in neural information processing systems*, 5, 1992.
- [35] Rotem Mulayoff, Tomer Michaeli, and Daniel Soudry. The implicit bias of minima stability: A view from function space. *Advances in Neural Information Processing Systems*, 34:17749–17761, 2021.
- [36] Nghia D Nguyen, Andrew Lutas, Jesseba Fernando, Josselyn Vergara, Justin McMahan, Jordane Dimidschstein, and Mark L Andermann. Cortical reactivations predict future sensory responses. *bioRxiv*, pages 2022–11, 2022.
- [37] Romain Brette. Is coding a relevant metaphor for the brain? *Behavioral and Brain Sciences*, 42:e215, 2019.
- [38] Elia Turner, Kabir V Dabholkar, and Omri Barak. Charting and navigating the space of solutions for recurrent neural networks. *Advances in Neural Information Processing Systems*, 34:25320–25333, 2021.

- [39] Honi Sanders, Matthew A Wilson, and Samuel J Gershman. Hippocampal remapping as hidden state inference. *Elife*, 9:e51140, 2020.
- [40] Blake A Richards, Timothy P Lillicrap, Philippe Beaudoin, Yoshua Bengio, Rafal Bogacz, Amelia Christensen, Claudia Clopath, Rui Ponte Costa, Archy de Berker, Surya Ganguli, et al. A deep learning framework for neuroscience. *Nature neuroscience*, 22(11):1761–1770, 2019.
- [41] Adam H Marblestone, Greg Wayne, and Konrad P Kording. Toward an integration of deep learning and neuroscience. *Frontiers in computational neuroscience*, page 94, 2016.
- [42] Andrew Saxe, Stephanie Nelli, and Christopher Summerfield. If deep learning is the answer, what is the question? *Nature Reviews Neuroscience*, 22(1):55–67, 2021.
- [43] Francisco M Codoñer, José-Antonio Darós, Ricard V Solé, and Santiago F Elena. The fittest versus the flattest: experimental confirmation of the quasispecies effect with subviral pathogens. *PLoS pathogens*, 2(12):e136, 2006.
- [44] Mor Shpigel Nacson, Rotem Mulayoff, Greg Ongie, Tomer Michaeli, and Daniel Soudry. The implicit bias of minima stability in multivariate shallow reLU networks. *inproceedings*, 2023.