From lazy to rich to exclusive task representations in neural networks and neural codes

Matthew Farrell, Stefano Recanatesi, and Eric Shea-Brown
March 2023

$_{\scriptscriptstyle 1}$ Abstract

- 2 Neural circuits both in the brain and in "artificial" neural network models learn to solve a remarkable
- ³ variety of tasks, and there is great current opportunity to use neural networks as models for brain function.
- ⁴ Key to this endeavor is the ability to characterize the *representations* formed by both artificial and biological
- $_{5}$ brains. Here, we investigate this potential through the lens of recently developing theory that characterizes
- $_{6}$ neural networks as "lazy" or "rich" depending on the approach they use to solve tasks: lazy networks
- $_{7}$ primarily solve tasks by selectively modifying readout weights, while rich networks solve tasks by modifying
- 8 weights throughout the network. We further elucidate rich networks through the lens of compression and
- 9 "neural collapse", ideas that have recently been of significant interest to neuroscience and machine learning.
- We then show how these ideas apply to a domain of increasing importance to both fields: extracting latent
- 11 structures through self-supervised learning.

12 Introduction

When we learn and develop, from learning to play chess to learning to walk to learning a relatively controlled laboratory task, the brain undergoes changes that specialize neural circuits to certain functions. However, the degree of specialization, and the elements of the environment and task that are specialized for, vary significantly across the brain. There is evidence that some brain areas hold a potentially high-dimensional (high-d), possibly random mix of many sensory features and task variables, as in high-d mixed-selective representations [2,3]. There is also evidence for other brain areas holding information that is very exclusively focused on task variables, such as category identity in classification tasks [4]. However, a unifying perspective that explains these phenomena is still developing, and it is not clear when and where general-purpose,

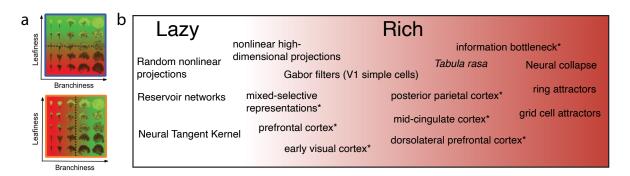


Figure 1: **a)** Depiction of a context-dependent classification task, taken from [1]. Top: task is to classify images based on leafiness. Bottom: task is to classify based on branchiness. **b)** An idealized schematic of the lazy to rich, and perhaps beyond, spectrum of representations that can follow from task learning. Asterisks denote items that are speculative or do not fit neatly into one place. "Prefrontal cortex" data are taken from [2] while other brain area data are taken from [1].

task-agnostic representations should be found versus exclusive, task-specific representations, and what the implications may be.

Artificial neural networks are proving a remarkably useful model system for addressing such questions [5]. In this approach, a neural network is trained to perform a task comparable to that performed by an animal in an experiment; the network is then analyzed to gain insight into plausible neural mechanisms and representations that support task execution. An obvious but appealing aspect of this approach is that representations in the artificial system can be measured with perfect completeness and precision. By studying the representations formed by neural networks, we can probe the functional role of representational structure. Moreover, by efficiently simulating and comparing representations that arise in many different settings, or performing mathematical analysis of well-defined learning processes, we can ask explicitly whether and how structure in these representations depends on the underlying tasks, and on varied assumptions about the networks and their learning rules. Of course, how much this informs the underlying biology is a deep, classic [6] and still open question, but a still-growing body of recent evidence underscores that the underlying representational principles may as well be at work in animal brains [7–9].

A striking finding that has recently emerged from theoretical investigations of neural networks is the large diversity of qualitatively different solutions these networks find, depending on initialization scheme, optimization procedure, and other details. These results belong to a sub-field that is often referred to as "feature learning." A feature is simply an aspect of inputs in a particular domain; feature learning describes how learning systems might access or extract useful features that support performing tasks within this domain. In neuroscience, the closely related concept of a neural representation is more familiar, which ultimately refers to how the biological network represents specific aspects of the external world. Often it is the case that useful features are encoded in neural representations. For instance, edges are useful features of images that support image classification; it is plausibly for this reason that edge-detecting simple cells are often found in primate visual cortex.

In this article, we highlight an intriguing set of recent findings about learned representations in artificial neural networks, and how they may shed light on biological neural representations and the underlying learning processes. In particular, we study:

- 1. When are networks *lazy*, in that they learn to accomplish tasks without changing their representations? When are networks instead *rich*, in that they change their representations over the course of learning?
- 2. In the case of rich networks, when does this richness reach the level of being *exclusive*, where their representations not just learn but also isolate task-relevant information?
- 3. How do exclusive networks shape representations in self-supervised tasks, such as predictive learning, to uncover hidden task structure or variables?

The third question has rapidly become of high importance to machine learning, as training regimens for artifical neural networks are increasingly dominated by a self-supervised initial stage, as well as in neuroscience, as self-supervised tasks are increasingly used as models to explain marquee neural representations such as place and grid cells.

Reviewing these topics will carry us across a spectrum from less to more extreme examples of feature learning: from no feature learning at all to exclusive feature learning that actively removes from representations any information inessential to the task at hand. Fig. 1a illustrates the concept of task-relevant and task-irrelevant features by showing a task where images of trees are classified either according to the tree leafiness or branchiness, according to a context signal. In the first context, aspects of the branchiness features are irrelevant to the leafiness features, and vice-versa in the second. Fig. 1b gives a rough schematic of the range of lazy to rich to exclusive learning by placing models, concepts, and brain areas along this spectrum.

Below, we begin with the regime in which features are not learned, known as the "lazy" regime in the feature learning literature. This regime has been important in the development of mathematical theories predicting the behavior of neural networks [10–13]. We then move to the regime in which feature learning occurs, known as the "rich" regime. Within the rich regime, we will investigate the different degrees of richness that can occur, which can be measured by the extent to which networks learn to exclusively represent information required for a specific task while rejecting other incoming information.¹ Throughout, we

¹Note that this terminology can be confusing, as the "rich" regime can involve *removing* information. Here richness refers to whether or not network representations are "enriched" (i.e. modified) over training in a way that reflects desired task outputs.

highlight insights that underlying concepts can provide for neurobiology (see, e.g., [1, 14, 15]).

72 The lazy learning limit: learning tasks without encoding them

Lazy learning describes a lack of task-relevant changes in the representation of a neural network (see below for a more formal definition). The most basic way in which lazy learning occurs is in model networks where internal connection weights are simply held constant by design: they are never allowed to vary from their initial values, which are chosen before learning begins. In this case, the only connection weights that are changed during learning are the output weights linking a network's internal activity to its final response. As a prominent example, networks whose internal connections are initialized with random internal connection weights may perform a sufficiently large number of transformations of their inputs in their internal layers (e.g. random nonlinear projections to a high-d space), so that the networks can solve many tasks by changes in readout weights alone. For example, in the illustration of Figure 2a, if the network is large, its neurons are nonlinear, and the internal (recurrent) weights W are sufficiently strong and random, then a vast set of input-output maps may be created by leaving W fixed and just tuning the readout weights. This is the approach taken by the related frameworks of support vector machines [16], random feature models [17], kernel machines [18], reservoir computers [19, 20], Koopman operators [21], and neural network Gaussian processes [22–25]. In particular, reservoir computers have featured prominently in the development of theories of neural computation [26]. This perspective of random nonlinear projections followed by learned readout weights is also prominent in circuit models for early sensory and cerebellar processing, e.g. [1,27-29].

Perhaps more surprisingly, even when we do allow synapses to be modified in both internal and readout layers (i.e., training the network "end-to-end"), lazy learning can still occur (cf. [11]). This phenomenon has come under intense focus in recent years. While networks will typically engage in some amount of task-relevant feature learning, it is theoretically convenient to consider limits in which feature learning either does or does not occur. One such limit is called the neural tangent kernel (NTK) limit. NTK theory says that if neural networks are initialized in a certain way, and if the size of each layer in the network is taken to infinity (as for neural circuits with vast numbers of internal neurons), then the evolution of the weights of the network through training – including the intermediate weights – can be described by a relatively simple mathematical object; namely, a linear system of ordinary differential equations. The coefficient matrix of this system is called the neural tangent kernel. This theory is enticing because the simple training dynamics allow theoreticians to analytically predict aspects of the behavior of the network through training, such as the error the network will have on held-out data after training a certain number of steps. These concepts have been applied to many network models, including convolutional neural networks as well as recurrent neural networks [31, 32]; however, simple unstructured feedforward networks are currently the best understood. Early evidence for the existence of a lazy regime for recurrent networks was also found in [33].

In the NTK limit, the network before training is randomly initialized, and the network weights after training (excepting the output weights) remain close to this random initialization, such that the intermediate representations are still essentially random. For this reason, networks initialized according to the Neural Tangent Kernel theory are said to be in the lazy regime. Due to this lack of feature learning, networks in the lazy regime resemble the support vector machines, random feature models, and other related ideas described above. A precise definition of the lazy regime can be found in [11]; for our purposes it is sufficient to understand that the intermediate feature representations do not change in a task-relevant way.² This said, many adjustments to the initialization scheme and other details can result in more substantial learning in internal representations [11]; in general much remains to be discovered (see [11–13,34–40]).

Interpretations for neurobiology

We pause to highlight three points important for relating the lazy regime to neurobiology. First, we reiterate that the scaling limits considered in these works are motivated mathematically, as they allow for a clean theory to be built, but that the ideas of rich and lazy can be used in a less formal way and applied to neural networks that are not infinite in size. In addition, while the formal theory as we have introduced it above

²To be a bit more precise, the training time would need to diverge to infinity along with the size of the network in order for features to be learned.

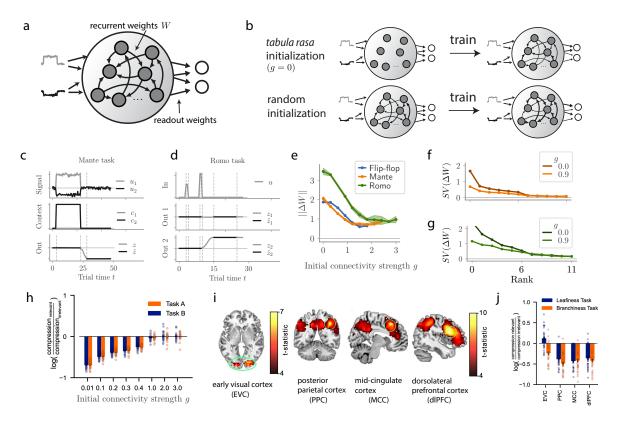


Figure 2: When do neural networks learn to represent task structure? Illustrating the influence of initial connectivity. a) Illustration of neural network. Inputs arrive and are transformed in internal, or "hidden," layer(s): here, this is the recurrent layer in the central circle. This internal representation is then read out, via readout weights, to become the output for the task. b) Cartoon comparing training from a tabula rasa initialization and a random initialization. Initial connectivity strength is denoted by q. Tabula rasa in this case corresponds to q=0. Panels c)-g) taken from [30]: c) Depiction of an instance of the "Mante task." Here, the network receives two noisy input signals, one of which is relevant in a given trial of the task and one of which serves as a distractor. The network also receives "context" inputs indicating the relevant signal. The network is tasked with outputting whether the relevant signal was positive or negative on average over the time course of the trial. d) Depiction of an instance of the "Romo task." Here, the network receives two input pulses separated by a delay, after which the network is tasked with outputting which of the two pulses had the larger amplitude. e) Frobenius norm of the change in weights resulting from training, as a function of initial connectivity strength. Colors denote different tasks. f) First 11 singular values of the change in weights resulting from training on the Mante task. Color denotes initial connectivity strength, g) As in f, but for the Romo task, h) - i) Adapted from [1], See [1] for details, h) Compression of task-irrelevant versus task-relevant information in a neural network with a single hidden layer trained on a simplified version of the context-modulated classification task illustrated in Fig. 1a, as a function of the initial connectivity strength q. More negative values indicate stronger compression of task-irrelevant information which is indicative of rich feature learning. i) Representation analysis of BOLD signals during human execution of the context-modulated classification task illustrated in Fig. 1a. Left panel: similarity of BOLD signals with input-specific features (signifying a lazy representation). Right three panels: similarity of BOLD signals with output/choice-specific features (signifying a rich representation). j) As in h, but measured for activity (representations) in different human brain regions during execution of the contextmodulated classification task illustrated in Fig. 1a.

refers to entire networks as potentially being lazy, we can use this concept for individual layers, or individual representations. Finally, we note that while lazy learning is defined above in terms of network initialization, there are multiple ways that this the concept of initialization could apply in biological settings; we discuss

118

119

120

Leaving room for learning: Rich learning of neural representations

Above, we reviewed how networks that start with strong internal network weights can show lazy learning, in which no meaningful learning of task structure occurs in a network's internal representations. The theoretical work underlying this phenomenon also suggests how networks may move away from this lazy learning as the (relative) strength of initial weights is tuned down (e.g., see [11]). Here we highlight some recent work from computational neuroscience that quantifies the type of learning that occurs in this setting.

28 Rich learning from scratch – the tabula rasa regime

The limit of taking very small initial weights is called the *tabula rasa* regime, and has played an important role in our understanding of network function. Here, the opposite of lazy behavior – termed *rich learning* – occurs: the task structure is strongly learned and represented in the underlying neural network. Fig. 2b shows an illustration which contrasts *tabula rasa* versus the strong and random initialization schemes similar to the NTK initialization discussed above. The behavior of the *tabula rasa* regime was made explicit in elegant mathematical studies of deep linear networks [35,41–43]. In this setting, the network picks up the structure of the task in a parsimonious fashion through training, with the modes of the input-output covariance matrix (the principal components) being transferred to the weights in order of their magnitude. In this way, network weights and activities clearly represent the structure of the task. With *tabula rasa* initialization, at least in the tractable case of linear feedforward networks, this is the only structure that weights represent.

139 Titrating away from tabula rasa

A recent study focused on recurrent neural networks explores the intermediate ground between larger and smaller norm initializations [30]. In this work, the *tabula rasa* behavior of the network is reconciled with a non-vanishing random initialization. This is done in the context of three tasks, two of which are shown in Fig. 2c and Fig. 2d. A key result is how random components of network weights present at initialization perturb the learning dynamics away from *tabula rasa* behavior. In particular, this random initialization is "sticky", with the network weights retaining higher-rank components through training (Figs. 2e to 2g). In this case the network after training may assume a lower-dimensional structure more dominated by a single mode if the random component is small or a higher-dimensional structure if the random component is large at initialization (Figs. 2f and 2g). In general, the main story that emerges seems to be that *changes* induced by learning have a rank that matches that of the task; often this is low-rank. See also [12, 35, 40, 44, 45] for explorations beyond the *tabula rasa* regime, which find similar principles at work.

In the cases explored in these studies, random components in the weights present at initialization tend to remain throughout training. In the section "Learning to be rich *and* exclusive" below, we will study the even more dramatic case of active compression of task-irrelevant information, where the random components present at initialization are significantly reduced through training.

155 Interpretations for neurobiology, revisited

We reiterate that both lazy and rich learning are defined above in terms of network initialization and how much the representation changes from this initialization. In biological circuits, what an initial network means is somewhat up to interpretation. For instance, a neural circuit may be considered *tabula rasa* at the outset of development (connections between neurons being weak or non-existent), and the learning process could be a mix of genetically-determined development along with synaptic modifications driven by experience in the world. Conversely, circuits that engage in learning a new tasks after a lifetime of learning other tasks may be considered to have stronger initial weights. Most settings will, of course, lie in between these extremes, and a great deal about the factors that control rich vs lazy learning outcomes doubtless remains to be discovered. In the meantime, the spectrum of representations that can arise in neural networks with the types of initializations studied to date form intriguing and testable hypotheses for experiments, as we review next.

167 Connections to experiments

Theoretical work exploring lazy and rich regimes has counterparts in experimental neuroscience studies such as [1]. In this study, the authors design a context-dependent task that can be solved in two ways: random projections to a high-d space followed by a trained readout (corresponding to the lazy regime), or using intermediate weights to transform the representation to a compact form (corresponding to the rich regime). In line with the computational work reviewed above, the authors found that a neural network trained on the task would take on the rich or lazy solution depending on the strength of connection weights in its (random) initialization (see Fig. 2h). They then studied neural activity during this task, and found that early sensory brain areas represent information in a way consistent with the lazy regime, whereas higher order areas such as posterior parietal cortex have representations more consistent with the rich regime (Figs. 2i and 2j).

Much work remains in charting out the strategies used in different brain areas across a range of tasks, and across different levels of functional hierarchy. To help address the need to understand brain representations across a large range of tasks, the authors of [46] compare across-task similarity matrices of brain representations to those formed by neural networks in the rich and lazy regimes, and found closer matches with rich networks. The study [47] also observes that deeper brain areas have a rich-learning-like representation, and investigates the generalization properties of these representations. These studies indicate the power and importance of the theory of rich and lazy networks when modeling neural representations.

Allied and very interesting concepts appear in the form of "mixed-selective" representations as observed in prefrontal cortex [2,3], which correspond to a high-d nonlinear mapping of sensory and task variables into the neural representation. We note that a mixed-selective representation is not necessarily equivalent to a lazy representation – the degree to which the representation is lazy depends on the degree to which the nonlinear mapping of input and task data is specifically tailored to the relevant tasks – a fully lazy representation will instead be random (see Fig. 1b for concrete examples). It is an interesting future direction to ascertain the degree to which mixed-selective and other neural representations are indeed random (see for example [48–50]), and to explore connections to mechanisms underlying lazy learning.

Learning to be rich and exclusive: compressing away task-irrelevant information

When networks operate outside the lazy regime, they change their internal representations in accordance with task demands. But do they become single-minded in this regard, compressing away input and information that is not directly relevant to the task at hand? This is a very strong way in which networks could encode tasks, since it isolates only the task-relevant information. It can involve the active removal of any information initially present but irrelevant to the task; signatures of this removal should be measurable in both experiments and model simulations. Such exclusive representations can also, in principle, have significant functional implications, such as enabling fast downstream learning and generalization on similar tasks [3,51–53], while limiting the ability of downstream networks to learn new tasks that require information that has been discarded [2,54].

Two complementary perspectives

We next review two perspectives on the compression of task-irrelevant information from the recent neural network literature. We first describe neural collapse, as it takes a direct "geometric" description of this compression that is easy to visualize. We then return to the earlier, and inspirational, idea of the information bottleneck, which quantifies compression of task-irrelevant inputs using mutual information.

Dimension compression and neural collapse

Recent studies [40,52,55–59] have taken a geometrical view on how networks can learn to actively compress away aspects of inputs that are not directly relevant to the task at hand, a phenomeon elegantly described by Papyan, Donoho, and colleagues as "neural collapse" [56]. Here, the structure of the recurrent neural network activity at later timesteps – or of deep neural network activity in penultimate layers – becomes very low-dimensional, in certain cases even collapsing to a set of single points. This occurs even though

the network is initialized with weights that form high-d representations carrying both task-relevant and task-irrelevant information.

In particular, for tasks with discrete outputs (or categories) into which inputs are grouped, the representation after training can become highly compressed around each category [52, 56, 59]. This behavior was observed in the context of change detection tasks [59] and discrete classification tasks [52], where it is referred to as dimensionality compression. Figures Fig. 3a-Fig. 3c, modified from [52], give an illustration of this phenomenon. Here, the task is delayed classification. First, an input arrives at timestep zero. A recurrent neural network then processes this input for some number of timesteps, until an evaluation time when the network activity is read out. Fig. 3b shows the phenomeon at hand: a network that has learned this task strongly compresses its inputs into clusters corresponding to the task categories. Here, the network at the initialization of learning approximately preserves the structure of the inputs; hence, arriving at a collapsed representation requires all of the task-irrelevant structure of the inputs to be quashed through learning [52]. Note that in this example, the two categories of inputs are linearly separable in the input space. Thus, the inputs do not need to be reformatted by the recurrent network in order to solve the task, as they could be classified with 100% accuracy by the output weights alone. Rather, the highly structured representations that form are an interesting by-product of how the network learns to solve the task. As in the preceding section, the initial connectivity strength plays a role in the nature of this compression; Fig. 3c shows how stronger initial connectivity leads to the formation of chaotic attractors that, while still compressed, are not as compressed as in the case of smaller initial connectivity strength (here both initializations are far from the tabula rasa regime).

Contemporaneously, a similar compression phenomenon was observed in convolutional neural networks trained on image recognition tasks [56] (Fig. 3d); see also [58,60]. The work of [56] also discovered further very interesting aspects of the geometry of the representation relevant to the higher-dimensional output space, such as the compressed clusters lying on the vertices of a simplex, and mathematically analyzed consequences of the resulting representations. We note that [61] highlighted a limitation of early studies of neural collapse, in that the phenomenon was examined for representations of training data and in cases may not be as robust for testing data.

Overall, we note that neural collapse adds additional structure beyond the low-rank representations that emerge in, e.g., linear networks trained in the *tabula rasa* regime [41, 42]. This is because the action of forming distinct localized clusters is highly nonlinear. This said, linearized analysis local to each separate cluster may still give insight into the underlying mechanisms [52, 58], especially if this linearization can be justified in some limit as is done in [45]. See also [40] for an in-depth analysis of compression in two-layer neural networks.

Information Bottleneck

The *information bottleneck* ideas and results by Tishby and colleagues preceded the papers above and played a highly influential role in showing how neural networks can compress task-irrelevant information overall [51,62]. The key quantity here is mutual information.

The authors demonstrate a very interesting phenomenon that can occur in deep neural networks: deep layers form representations which learn to maximize mutual information about task outputs, while minimizing overall information about network inputs. This is illustrated via the "learning curves" in Fig. 3g, from [51,63]. For a given layer T, these curves track two quantities over the course of network learning. The first, I(X,T), is the mutual information between representations in layer T and the network inputs, X. The second, I(T,Y), is the corresponding information for the task outputs, Y. A key point is that, as learning evolves across epochs and I(T,Y) continues to increase, eventually I(X,T) begins to decrease. Thus, task-irrelevant mutual information is gradually compressed over the course of learning.

Remarkably, the authors also show that this compression can follow a precise optimization relationship, in which the mutual information about task inputs is minimized, subject to preserving mutual information about task outputs. This connects directly to earlier analytical work on the broader concept of an information bottleneck [64]. Moreover, the authors develop a mechanism by which this bottleneck may develop through the course of learning in neural netowrks. This stochasticity is inherent in incremental learning processes, in which network weights are incremented step by step based the successive examples (or "batches" of examples) on which the network is trained. Because each is randomly selected, there is a deviation on each step from

the true task gradient. Intriguingly, the authors show how this can, at least under certain assumptions, result in the same selective compression of task-irrelevant input information that leads to the information bottleneck found in their simulations and broader theory.

Robustness of exclusive representations

As in the results reviewed above showing that the emergence of rich vs. lazy representations is not depends (at least) on network initialization, the emergence of representations that compress task-irrelevant information is also far from automatic. Rather, it depends on many variables including network architecture, initialization scheme, loss function, and optimization procedure. In fact, for some tasks that lead to highly compressed representations in some settings, even changes to the loss function and details of gradient learning algorithm ("optimizer"), can eliminate compression. As a striking example, the compression seen in Figs. 3a and 3b vanishes when the optimizer is changed from RMSprop to "vanilla" stochastic gradient descent [52], and the network instead uses a lazy strategy to solve the task.

Given this possibility, when does neural collapse actually occur? The study [52] highlights three mechanisms that encourage neural collapse: the combination of loss functions that encourage scaling up outputs (such as categorical cross entropy) with saturating nonlinearities in the intermediate layers, large variability during the training process, and weight decay (in which additional terms are added to the network cost function to penalize large weights during training). The underlying theory generally uses a geometric decomposition of the learning dynamics into task-relevant and task-irrelevant directions [52,57,58]. In Figs. 3e and 3f, we provide supporting evidence for this theory in the context of the deep convolutional neural network Resnet18 [65] trained on the image recognition task CIFAR-10 [66]. Here, we find that variability induced by dropout – where neurons are randomly silenced, similar to a decrease in firing rate caused by a biological neuron's failure to spike – results in more compressed representations (Fig. 3e for the RMSprop optimizer, and Fig. 3f for stochastic gradient descent with momentum). This suggests that noise inherent in biological neural circuits may drive neural collapse. We use a simple measure of task-irrelevant compression, which is simply the average within-class distance divided by the average across-class distance of points in representation space. This collapse is measured on testing data, with mean squared error as a loss function (compare with [61] where categorical cross-entropy loss is used).

Similarly, not all networks, tasks, and training processes lead to information bottlenecks. For example, [63] argue that this result can be limited to networks with double-sided saturating nonlinearities such as the hyperbolic tangent functions, again limiting the universality of the phenomenon (Fig. 3h).

Titrating the representation: partial compression across time and across layers

A highly interesting aspect of the compression of task-irrelevant information in neural networks is that, even in fully trained networks, it is not necessarily an all or none process. Rather, several studies have shown that this compression can occur gradually across layers of deep networks, or gradually across time in recurrent neural networks. This gradation has potentially important implications in neuroscience: downstream brain areas with access to the neural network at different stages (layers or timesteps) would have access to different levels of stimulus information. This could be useful in driving different learning, memory, or behavioral systems downstream ([67] and John Maunsell, personal communication).

The graded nature of information compression is clearly evident in the information bottleneck formulation, where learning occurs in a two-phase process. In the first stage, intermediate layers first gain information about inputs and outputs (with input information dominating in the earlier layers and output information dominating in later layers). In the second phase, information about inputs is progressively lost (at a rate that is faster for later layers) [51,62] (Fig. 3g). The authors of [52,58,60] also study collapse over layers, or timesteps in the case of recurrent neural networks. They find that compression occurs progressively over the timesteps of the recurrent network (or over layers of a feedforward network). However, when trained on low-dimensional input data, the network first lifts the representation into a higher-dimensional space in the first few timesteps, reminiscent of the kernel machine approach, before potentially compressing the representation back down in later timesteps ([52], see also [68]).

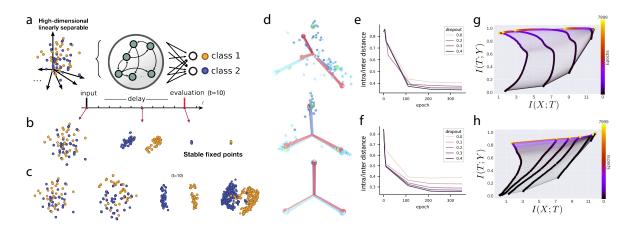


Figure 3: When do neural networks compress task-irrelevant inputs and information? Illustrating neural collapse and the information bottleneck. a) Illustration of a recurrent network solving a delayed classification task with high-d, and hence linearly separable, inputs; see text. The network's cost function is defined via categorical cross-entropy, with no additional penalty terms added. b) Visualization of the trained network representation's top two principal components as it evolves through time. c) As in b, but for stronger initial connectivity strength. d) Illustration of the representation of the penultimate layer of a convolutional neural network (adapted from [56]), projected onto the top 3 principle component axes (axes not shown). Green spheres denote axes of a 2-dimensional simplex, red balls and sticks represent the normal vectors for classifying hyperplanes, blue balls and sticks denote class means, and small blue spheres represent last-layer features. Top is before training, middle is at an intermediate stage of training, and bottom is after training. Note that this representation is induced by the training set, not the test set. e)-f) Plots of average intra-class divided by inter-class distances over training in the penultimate layer of Resnet18, where dropout is applied during training (but not during evaluation of distances). Lower values indicate more task-irrelevant compression. Loss used is mean squared error. Shading denotes amount of dropout applied. e) Optimizer used is stochastic gradient descent with momentum value of 0.9, a common optimizer choice in training neural networks. f) Optimizer is RMSprop without momentum, another common choice. g) Mutual information between a convolutional neural network layer's representation T and the inputs X(x-axis) as well as the outputs Y (y-axis) in a convolutional neural network with tanh nonlinearities (taken from [51,63]). Each curve corresponds with a different layer T and color denotes training time (epochs). h) As in **g**, but with ReLU nonlinearities (taken from [63]).

Connections to experiments

314

315

316

317

318

319

321

322

323

324

325

326

327

328

320

330

Many studies have observed neural representations that are highly compressed when compared to nonlinear high-d representations. These include [1, 4, 47], which show that particular brain areas such as posterior parietal cortex are relatively compressed (see Figs. 2h to 2j). In addition, [4] shows that lateral intraparietal (LIP) neurons form a highly-compressed representation that reflects the structure of task outputs, and that adapts to new task structures. In contrast, the middle temporal areas have a lazy representation that reflects the structure of the inputs and does not adapt to changing task structure

Less is known about whether certain brain areas/representations truly compress away information that is irrelevant to the task at hand, or at what temporal scale information is lost. To quote [4], "The exact nature of the role of LIP during learning, and whether changes in the ... representations of LIP are stable or vary dynamically with the demands of the task, remain to be determined." While these questions still await more final answers, they can be approached by experimental paradigms that track learning over time. Indeed, the work of Stern and colleagues [59] both identified compression of representation dimension computationally and found evidence for this compression gradually emerging over days of task learning in widefield activity patterns of mice. The degree to which the process of compressing task-irrelevant information – or "learning to forget" – can be tracked dynamically in the brain has important implications for our understanding of brain function, and remains an exciting direction for further experimental studies.

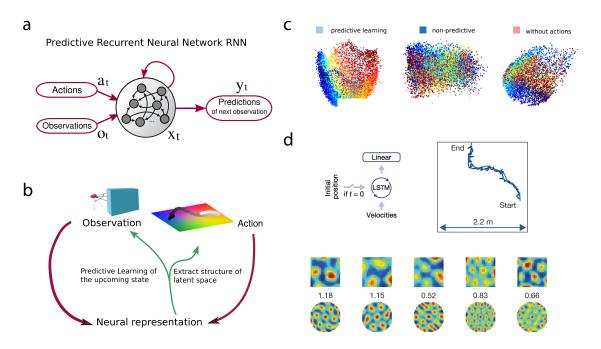


Figure 4: Learning to represent information in the absence of labeled data: Predictive learning in navigation tasks and emergent latent variables. a-c adapted from [69]. a) Illustration of the network architecture for a navigational predictive learning task. b) Illustration of the task environment for a navigational predictive learning task. c) The top three principal components for the representation of a neural network after training the predictive task, colored by position and orientation, and compared with other versions of the task. Left: predictive learning. Center: non-predictive learning. Right: predictive learning without the action information. d) Emergence of grid cells in a recurrent neural network trained to path integrate. Note that path integration is a form of prediction – velocity signals (actions) are used to predict future positions. Taken from [70] (reproduced with permission from Nature).

Exclusive network representations in self-supervised tasks

Above, we have focused on supervised learning tasks, in which every input is paired with an intended output – for example images and their labels. However, a vast and rapidly growing field studies neural networks trained on data without labels beyond that provided by the inputs themselves – a class of tasks known as semi-supervised or self-supervised (see, e.g., [71–74]). Importantly, this setting likely corresponds to much of the learning that occurs in neurobiology as well, as in the real world there is usually no explicit "teacher" constantly defining correct task outputs but rather sustained exposure to a richly informative sensory world.

Rather than attempting to cover the vast literature on semi-supervised learning and representations, we ask the specific question that ties most immediately to the above: what are the consequences of rich and exclusive representations for self-supervised learning? We leave the possibility of lazy solutions to predictive learning as an interesting topic for further exploration.

We focus on the prominent self-supervised setting of predictive learning. In predictive learning inputs occur as a time series and networks are trained to predict their future values. Here, the target output of the network is a temporally-shifted version of the input, which the network is trained to predict. Thus, task-relevant features are features that bear predictive power over future inputs. Conversely, task-irrelevant features do not influence future inputs.

Collapsed representations and the extraction of latent variables

A prominent example of a predictive learning task is navigation. In one recent study [69], a recurrent network was trained to predict future visual observations of an agent moving through an environment (Figs. 4a

and 4b). Here, the agent's actions were part of the input but not the output of the network, thus reflecting the agent's internal knowledge. In such a task the agent's field of view, and thus its visual observations, are completely determined by its state in the environment, the (x,y) location and orientation θ , so that (x,y,θ) are the underlying latent variables or latent states. This means that if the network could infer the current latent state (x,y,θ) from the visual inputs, it could in principle apply the agent's actions and anticipate future observations, through a downstream mapping from the latent space to observations. The work [69] showed that predictive learning extracts such latent states, through activity patterns that emerge in a trained recurrent neural network. Moreover, such latent states are less apparent in control networks trained to encode, but not predict, their inputs (Fig. 4c).

In this sense, the predictive network learns a parsimonious model of the structure of its world, and rules for how to update the states defining that world based on its actions. In the process, detailed visual scene information is compressed. This results in relatively low-dimensional neural representations of the latent states. Recall that the task-relevant features here are those that enable the prediction of future inputs – that is, the latent states. Thus, the network has learned an internal representation that is compressed around these task-relevant features.

The above gives a concrete illustration of how the ideas of representing task structure and compressing task-irrelevant inputs can be applied to the setting of predictive learning: the sole distinction from the supervised cases above is that the concept of task-relevant structure has been replaced by the concept of latent-state structure [69]. Recently, biologically-inspired Hebbian learning rules have also been derived that perform predictive learning [75, 76], and the network behavior (and extracted latent variables) analyzed.

In theory, networks in the lazy regime will not form structured representations such as these. As far as we are aware, whether or not networks in the lazy regime are able to satisfactorily perform predictive learning remains an open avenue for future work.

Connections to neuroscience

Several influential lines of work propose that predictive learning is a major driver of neural representations across the brain [77–79]. This said, the modeling approaches proposed for different brain areas can differ. Models focusing on the visual stream (i.e. sensory prediction) have historically been decoupled from those focusing on hippocampal dynamics (i.e. memory-based prediction). Moreover, the former models have placed emphasis on spatial prediction (e.g. completing missing elements of an image), and the latter on temporal prediction (similar to that discussed above). In the first instance, predictive models have been found to extract representations that reproduce visual receptive fields [80] and other properties of how sensory systems encode information [75,81]. In the second, the majority of research has centered on navigation, demonstrating that during predictive learning neurons begin to tile location and orientation in their activations [69]. This is similar to place cells and head direction cells in navigation-related brain circuits [82].

The case of grid cells merits special consideration. In [70], the authors trained a network to directly predict the next (x, y) location based on the current location, orientation and the action the agent takes (a calculation known as path integration) (Fig. 4d). The authors found that many of the units in the trained network functioned as grid cells, thus adopting a well-known encoding of spatial latent variables. While a theoretical accounting of this phenomenon was provided by [83], the sensitivity of grid cells' emergence to choice network architectures, and training rules, and allied hyperparameters is a topic of ongoing research [70, 84, 85]. As one example, the authors of [70] note dropout as an important mechanism in the appearance of grid cell representations, perhaps connecting to work on the principles of neural collapse reviewed above.

Discussion

The prevalalence of task-trained neural networks as models of the brain is exploding. This makes understanding the robustness and universality of neural network behavior essential, so that we can properly contextualize the insights these behaviors may provide for neural circuits in the brain. In this spirit, we began our review with a fundamental but remarkably subtle question: When do such networks learn to encode tasks in an observable way? First, we introduced multiple regimes to determine when neural networks learn task features in their internal representations. Then, we investigated the allied phenomenon of

forming compressed representations that isolate information that is relevant to the task at hand. Finally, we demonstrated how these principles arise in semi-supervised learning.

The emergence of task-relevant features in neural representations is not a given

While both feature learning and compressed/exclusive representations occur widely, they do not occur all the time. This is a major takeaway from our review: the details of learning algorithms and network structure matter. It may not be sufficient to train a neural network and examine the learnt solutions; rather, it is necessary to verify, and illuminating to explore, the robustness of the findings about network representations with respect to a number of factors, including initialization and optimization schemes. One extreme example is that neural networks can sometimes learn to solve tasks by merely updating output weights during training, or by updating weights in a manner that does not truly lead to learning task-driven features (such as with NTK initialization). At another extreme, they may begin with no structure and learn task-driven features alone (tabula rasa), or begin with task-irrelevant structure that is removed through training (compression/neural collapse). Intermediate outcomes are also possible.

As such, training a single network as a model of a brain circuit is unlikely to bear definitive results. Rather, we will need to evaluate the spectrum of responses that neural networks can exhibit, or clearly justify specific choices of their initialization, architecture, and learning rules. Theory can help guide the way here. As we reviewed, some of the mechanisms that appear to encourage task-structured representations include weight decay, "noisy" optimization processes such as RMSprop and ADAM, saturating nonlinearities, and added sources of network noise such as dropout. This said, much remains to be understood about the role of these and other network and learning mechanisms.

Reinforcement learning, and the revenge of the "hand-built" model

While we have focused this review on the behavior of trained neural networks, there are many other domains in which the distinction between lazy and rich solutions is highly relevant. One example are networks trained via reinforcement learning algorithms. It is still very much an open challenge to delineate the two regimes in this domain, and to work out the functional consequences of the different approaches.

Another important example are networks where many components are not trained, but "hand-built" with a given functionality in mind. In this approach, typically a theorist has in mind a desired network function, and uses geometrical and mathematical reasoning to build networks that fullfill it. In some instances only the output weights are trained, and in others there is no need for training at all. An example of the latter are attractor networks hand-built to integrate velocity signals in order to track position, such as models of head-direction cells [82] and grid cells [86]. Indeed, linking network models of brain circuits provided by gradient descent training procedures and those that could be, at least in principle, built by hand, provides a fruitful path toward interpreting and understanding how the trained networks actually work (cf. [87]). While hand-built models tend to resemble networks in the rich regime, some models, such as reservoir computers, are designed to be in the lazy regime.

Predictive learning in language models and beyond

We end our review with a discussion on the burgeoning area of self-supervised predictive learning. Prediction forces a learning agent to learn how operations and actions influence the world. Such an understanding is likely best supported by an efficient representation that reflects the relatively low-dimensional structure of the latent variables, though this remains to be proved. Here, we discussed likely connections to feature learning and information compression. However, theory describing the structure (or lack thereof) extracted by predictive learning, especially under the wide possibilities of different initializations, optimization schemes, and other hyperparameters, is still in its infancy.

Language models are a particularly topical example of the ability of predictive models to extract latent space information. While we make no attempt to review the tremendous advances made in natural language processing (NLP) over the past decade, we note that the underlying models have employed multiple prediction-based techniques to extract language structure. For instance, the famous project word2vec [88] reveals that neural networks trained to predict omitted words acquire a representation of these words that

forms a latent space map corresponding to word meanings. With this representation, word-based manipulations appear meaningful (e.g. King - Man + Woman = Queen). Recent theoretical work has enhanced our understanding of such vector-symbolic operations [89–91]. The efficacy of allied models [92] in extracting the fundamental structure of languages has led to breakthroughs such as ChatGPT and crosslingual translation between any two languages [93–95].

This said, while predictive training has been a cornerstone in developing language models, experimental tests of whether the underlying representations also appear in biological brains have been limited by the difficulty of conducting language-based studies in nonhuman animals. Nevertheless, recent research indicates that human-level language comprehension involves predictive processing [96, 97]. Finally, we note that implicit in the above is our speculation that the emergence of latent variable structure in the representation of language is driven, at least in part, by the same factors reviewed above that promote the compression of task-irrelevant information. Verifying or rejecting this speculation is an intriguing target for future modeling and theoretical work.

460 Acknowledgements

We thank the reviewers for their careful and thoughtful suggestions which have substantially improved this manuscript. E. S-B. gratefully acknowledges support of NIH Grants UF1NS126485 and RF1DA055669 and of NSF NCS-FO Grant 2024364. All authors gratefully thank the Swartz Foundation and its Centers at Harvard (MF) and at the University of Washington for their support.

465 References

447

448

449

451

452

453

455

456

457

458

459

- Timo Flesch, Keno Juechems, Tsvetomira Dumbalska, Andrew Saxe, and Christopher Summerfield. Orthogonal representations for robust context-dependent task performance in brains and neural networks.
 Neuron, 110(7):1258–1270.e11, April 2022.
- Mattia Rigotti, Omri Barak, Melissa R. Warden, Xiao-Jing Wang, Nathaniel D. Daw, Earl K. Miller, and
 Stefano Fusi. The importance of mixed selectivity in complex cognitive tasks. *Nature*, 497(7451):585–590, May 2013.
- [3] Stefano Fusi, Earl K Miller, and Mattia Rigotti. Why neurons mix: High dimensionality for higher cognition. Current Opinion in Neurobiology, 37:66–74, April 2016.
- ⁴⁷⁴ [4] David J. Freedman and John A. Assad. Experience-dependent representation of visual categories in parietal cortex. *Nature*, 443(7107):85–88, September 2006.
- ⁴⁷⁶ [5] Omri Barak. Recurrent neural networks as versatile tools of neuroscience research. Current opinion in neurobiology, 46:1–6, 2017.
- [6] David E. Rumelhart and James L. McClelland. Parallel Distributed Processing: Explorations in the Microstructure of Cognition: Foundations. MIT Press, 1987.
- [7] Guangyu Robert Yang and Xiao-Jing Wang. Artificial neural networks for neuroscientists: a primer.
 Neuron, 107(6):1048–1070, 2020.
- ⁴⁸² [8] Valerio Mante, David Sussillo, Krishna V. Shenoy, and William T. Newsome. Context-dependent computation by recurrent dynamics in prefrontal cortex. *Nature*, 503(7474):78–84, November 2013.
- [9] Daniel L. K. Yamins, Ha Hong, Charles F. Cadieu, Ethan A. Solomon, Darren Seibert, and James J.
 DiCarlo. Performance-optimized hierarchical models predict neural responses in higher visual cortex.
 Proceedings of the National Academy of Sciences, 111(23):8619–8624, 2014.
- [10] Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural Tangent Kernel: Convergence and Generalization in Neural Networks. Advances in Neural Information Processing Systems, page 10, 2018.

- [11] Greg Yang and Edward J. Hu. Tensor programs iv: Feature learning in infinite-width neural networks.

 In International Conference on Machine Learning, 2021.
- [12] Jacob A Zavatone-Veth, Abdulkadir Canatar, Benjamin S Ruben, and Cengiz Pehlevan. Asymptotics
 of representation learning in finite Bayesian neural networks. Journal of Statistical Mechanics: Theory
 and Experiment, 2022(11):114008, November 2022.
- ⁴⁹⁴ [13] Blake Bordelon and Cengiz Pehlevan. The influence of learning rule on representation dynamics in wide neural networks. In *The Eleventh International Conference on Learning Representations*, 2023.
- ⁴⁹⁶ [14] Haitao Zhao, Zhihui Lai, Henry Leung, and Xianyi Zhang. Feature learning and understanding: Algorithms and Applications. Springer Cham, 2020.
- [15] SueYeon Chung and LF Abbott. Neural population geometry: An approach for understanding biological and artificial neural networks. *Current opinion in neurobiology*, 70:137–144, 2021.
- [16] Ingo Steinwart and Andreas Christmann. Support Vector Machines. Information Science and Statistics.
 Springer, New York, 1st ed edition, 2008.
- [17] Ali Rahimi and Benjamin Recht. Uniform approximation of functions with random bases. In 2008 46th Annual Allerton Conference on Communication, Control, and Computing, pages 555–561, 2008.
- [18] Bernhard Schölkopf and Alexander J. Smola. Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond. Adaptive Computation and Machine Learning. MIT Press, 2002.
- 507 [19] Wolfgang Maass, Thomas Natschläger, and Henry Markram. Real-Time Computing Without Stable
 508 States: A New Framework for Neural Computation Based on Perturbations. Neural Computation,
 509 14(11):2531–2560, November 2002.
- [20] Herbert Jaeger. The "echo state" approach to analysing and training recurrent neural networks-with an
 erratum note. Bonn, Germany: German National Research Center for Information Technology GMD
 Technical Report, 148:1–47, January 2001.
- [21] Igor Mezić. Koopman Operator, Geometry, and Learning of Dynamical Systems. *Notices of the American Mathematical Society*, 68(07):1, August 2021.
- [22] Christopher Williams. Computing with infinite networks. In M.C. Mozer, M. Jordan, and T. Petsche, editors, *Advances in Neural Information Processing Systems*, volume 9. MIT Press, 1996.
- 517 [23] Radford M. Neal. Priors for Infinite Networks, pages 29–53. Springer New York, New York, NY, 1996.
- Jaehoon Lee, Jascha Sohl-dickstein, Jeffrey Pennington, Roman Novak, Sam Schoenholz, and Yasaman Bahri. Deep neural networks as gaussian processes. In *International Conference on Learning Representations*, 2018.
- [25] Alexander G. de G. Matthews, Jiri Hron, Mark Rowland, Richard E. Turner, and Zoubin Ghahramani.

 Gaussian process behaviour in wide deep neural networks. In *International Conference on Learning Representations*, 2018.
- [26] Steven Chase, Andrew Schwartz, Wolfgang Maass, and Robert Legenstein. Functional network reorganization in motor cortex can be explained by reward-modulated Hebbian learning. In Y. Bengio,
 D. Schuurmans, J. Lafferty, C. Williams, and A. Culotta, editors, Advances in Neural Information Processing Systems, volume 22, pages 1105–1113. Curran Associates, Inc., 2009.
- ⁵²⁸ [27] Ashok Litwin-Kumar, Kameron Decker Harris, Richard Axel, Haim Sompolinsky, and L. F. Abbott. Optimal Degrees of Synaptic Connectivity. *Neuron*, 93(5):1153–1164.e7, March 2017.
- [28] Frederic Lanore, N. Alex Cayco-Gajic, Harsha Gurnani, Diccon Coyle, and R. Angus Silver. Cerebellar granule cell axons support high-dimensional representations. *Nature Neuroscience*, June 2021.

- [29] N. Alex Cayco-Gajic, Claudia Clopath, and R. Angus Silver. Sparse synaptic connectivity is required for decorrelation and pattern separation in feedforward networks. *Nature Communications*, 8(1):1–11, October 2017.
- Friedrich Schuessler, Francesca Mastrogiuseppe, Alexis Dubreuil, Srdjan Ostojic, and Omri Barak. The
 interplay between randomness and structure during learning in rnns. In H. Larochelle, M. Ranzato,
 R. Hadsell, M.F. Balcan, and H. Lin, editors, Advances in Neural Information Processing Systems,
 volume 33, pages 13352–13362. Curran Associates, Inc., 2020.
- [31] Sina Alemohammad, Zichao Wang, Randall Balestriero, and Richard Baraniuk. The recurrent neural tangent kernel. In *International Conference on Learning Representations*, 2021.
- [32] Greg Yang and Etai Littwin. Tensor programs iib: Architectural universality of neural tangent kernel
 training dynamics. In Marina Meila and Tong Zhang, editors, Proceedings of the 38th International
 Conference on Machine Learning, volume 139 of Proceedings of Machine Learning Research, pages
 11762–11772. PMLR, 18–24 Jul 2021.
- [33] Ulf D. Schiller and Jochen J. Steil. Analyzing the weight dynamics of recurrent learning algorithms.
 Neurocomputing, 63:5–23, 2005.
- Andrew Saxe, Shagun Sodhani, and Sam Jay Lewallen. The neural race reduction: Dynamics of abstraction in gated networks. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang
 Niu, and Sivan Sabato, editors, Proceedings of the 39th International Conference on Machine Learning,
 volume 162 of Proceedings of Machine Learning Research, pages 19287–19309. PMLR, 17–23 Jul 2022.
- [35] Alexander Atanasov, Blake Bordelon, and Cengiz Pehlevan. Neural networks as kernel learners: The silent alignment effect. In *International Conference on Learning Representations*, 2022.
- [36] Ethan Dyer and Guy Gur-Ari. Asymptotics of wide networks from feynman diagrams. In *International Conference on Learning Representations*, 2020.
- Qianyi Li and Haim Sompolinsky. Statistical mechanics of deep linear neural networks: The backpropagating kernel renormalization. *Phys. Rev. X*, 11:031059, Sep 2021.
- [38] Jacob A. Zavatone-Veth, William L. Tong, and Cengiz Pehlevan. Contrasting random and learned
 features in deep bayesian linear regression. Phys. Rev. E, 105:064118, Jun 2022.
- [39] Daniel A. Roberts, Sho Yaida, and Boris Hanin. The Principles of Deep Learning Theory. Cambridge
 University Press, 2022. https://deeplearningtheory.com.
- [40] Jonas Paccolat, Leonardo Petrini, Mario Geiger, Kevin Tyloo, and Matthieu Wyart. Geometric compression of invariant manifolds in neural networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2021(4):044001, apr 2021.
- [41] Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics
 of learning in deep linear neural networks. *International Conference on Learning Representations 2014*,
 February 2014.
- [42] Andrew M. Saxe, James L. McClelland, and Surya Ganguli. A mathematical theory of semantic development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116(23):11537–11546,
 June 2019.
- [43] Jianghong Shi, Eric Todd SheaBrown, and Michael A Buice. Learning dynamics of deep linear networks
 with multiple pathways. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho,
 editors, Advances in Neural Information Processing Systems, 2022.
- Lukas Braun, Clémentine Carla Juliette Dominé, James E Fitzgerald, and Andrew M Saxe. Exact learning dynamics of deep linear networks with prior knowledge. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, Advances in Neural Information Processing Systems, 2022.

- Francesca Mastrogiuseppe, Naoki Hiratani, and Peter Latham. Evolution of neural activity in circuits bridging sensory and abstract knowledge. *eLife*, 12:e79908, March 2023.
- Takuya Ito and John D. Murray. Multitask representations in the human cortex transform along a sensory-to-motor hierarchy. *Nature Neuroscience*, 26(2):306–315, February 2023.
- [47] W. Jeffrey Johnston and Stefano Fusi. Abstract representations emerge naturally in neural networks trained to perform multiple tasks. *Nature Communications*, 14(1):1040, February 2023.
- Alexis Dubreuil, Adrian Valente, Manuel Beiran, Francesca Mastrogiuseppe, and Srdjan Ostojic. The role of population structure in computations through neural dynamics. *Nature Neuroscience*, 25(6):783–794, June 2022.
- Junya Hirokawa, Alexander Vaughan, Paul Masset, Torben Ott, and Adam Kepecs. Frontal cortex
 neuron types categorically encode single decision variables. Nature, 576(7787):446-451, December 2019.
- David Raposo, Matthew T. Kaufman, and Anne K. Churchland. A category-free neural population supports evolving demands during decision-making. *Nature neuroscience*, 17(12):1784–1792, December 2014.
- [51] Naftali Tishby. The Information Bottleneck Theory of Deep Neural Networks. Bulletin of the American
 Physical Society, 2018.
- [52] Matthew Farrell, Stefano Recanatesi, Timothy Moore, Guillaume Lajoie, and Eric Shea-Brown.
 Gradient-based learning drives robust representations in recurrent neural networks by balancing compression and expansion. Nature Machine Intelligence, 4:1–10, June 2022.
- [53] Abdulkadir Canatar, Blake Bordelon, and Cengiz Pehlevan. Spectral bias and task-model alignment
 explain generalization in kernel regression and infinitely wide neural networks. Nature Communications,
 12(1):2914, May 2021.
- [54] Silvia Bernardi, Marcus K. Benna, Mattia Rigotti, Jérôme Munuera, Stefano Fusi, and C. Daniel Salzman. The geometry of abstraction in the hippocampus and prefrontal cortex. Cell, 183(4):954–967.e21, 2020.
- [55] Matthew Farrell. Revealing structure in trained neural networks through dimensionality-based methods.

 *University of Washington ProQuest Dissertations Publishing, 2020.**
- [56] Vardan Papyan, X. Y. Han, and David L. Donoho. Prevalence of neural collapse during the terminal
 phase of deep learning training. Proceedings of the National Academy of Sciences, 117(40):24652–24663,
 October 2020.
- [57] X.Y. Han, Vardan Papyan, and David L. Donoho. Neural collapse under MSE loss: Proximity to and dynamics on the central path. In *International Conference on Learning Representations*, 2022.
- [58] Stefano Recanatesi, Matthew Farrell, Madhu Advani, Timothy Moore, Guillaume Lajoie, and Eric Shea Brown. Dimensionality compression and expansion in Deep Neural Networks. arXiv:1906.00443 [cs, stat], October 2019.
- [59] M. et al. Stern. In the footsteps of learning: Changes in network dynamics and dimensionality with task acquisition. COSYNE Conference Abstract, 2020.
- [60] Alessio Ansuini, Alessandro Laio, Jakob H Macke, and Davide Zoccolan. Intrinsic dimension of data
 representations in deep neural networks. NIPS, page 11, 2019.
- [61] Like Hui, Mikhail Belkin, and Preetum Nakkiran. Limitations of neural collapse for understanding generalization in deep learning. arXiv:2202.08384v1 [cs.LG], 2022.
- [62] Ravid Shwartz-Ziv and Naftali Tishby. Opening the Black Box of Deep Neural Networks via Information.

 arXiv:1703.00810v3 [cs.LG], March 2017.

- [63] Andrew Michael Saxe, Yamini Bansal, Joel Dapello, Madhu Advani, Artemy Kolchinsky, Brendan Daniel
 Tracey, and David Daniel Cox. On the information bottleneck theory of deep learning. In *International Conference on Learning Representations*, 2018.
- [64] Naftali Tishby, Cicero Pereira, and William Bialek. The information bottleneck method. Proceedings
 of the 37th Allerton Conference on Communication, Control and Computation, 49, 07 2001.
- [65] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 770–778, 2016.
- [66] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- [67] Cory Stephenson, suchismita padhy, Abhinav Ganesh, Yue Hui, Hanlin Tang, and SueYeon Chung. On
 the geometry of generalization and memorization in deep neural networks. In *International Conference* on Learning Representations, 2021.
- [68] Christian Keup, Tobias Kühn, David Dahmen, and Moritz Helias. Transient chaotic dimensionality
 expansion by recurrent networks. Phys. Rev. X, 11:021064, Jun 2021.
- [69] 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):1417, March 2021.
- [70] Andrea Banino, Caswell Barry, Benigno Uria, 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, 2018.
- [71] Longlong Jing and Yingli Tian. Self-supervised visual feature learning with deep neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 43(11):4037–4058, 2020.
- [72] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation Learning with Contrastive Predictive
 Coding. arXiv:1807.03748 [cs, stat], January 2019.
- [73] Geoffrey Hinton and Terrence J. Sejnowski. *Unsupervised Learning: Foundations of Neural Computa*tion. The MIT Press, May 1999.
- Randall Balestriero, Mark Ibrahim, Vlad Sobal, Ari Morcos, Shashank Shekhar, Tom Goldstein, Florian
 Bordes, Adrien Bardes, Gregoire Mialon, Yuandong Tian, Avi Schwarzschild, Andrew Gordon Wilson,
 Jonas Geiping, Quentin Garrido, Pierre Fernandez, Amir Bar, Hamed Pirsiavash, Yann LeCun, and
 Micah Goldblum. A Cookbook of Self-Supervised Learning. arXiv:2304.12210 [cs], June 2023.
- [75] David Lipshutz, Charles Windolf, Siavash Golkar, and Dmitri Chklovskii. A biologically plausible neural network for slow feature analysis. Advances in neural information processing systems, 33:14986–14996,
 2020.
- [76] Manu Srinath Halvagal and Friedemann Zenke. The combination of hebbian and predictive plasticity
 learns invariant object representations in deep sensory networks. bioRxiv, 2023.
- ⁶⁵⁸ [77] Christoph Teufel and Paul C Fletcher. Forms of prediction in the nervous system. *Nature Reviews* Neuroscience, 21(4):231–242, 2020.
- [78] Howard Eichenbaum and Norbert J Fortin. The neurobiology of memory based predictions. Philosophical
 Transactions of the Royal Society B: Biological Sciences, 364(1521):1183–1191, 2009.
- [79] Rajesh PN Rao. A sensory-motor theory of the neocortex based on active predictive coding. *bioRxiv*, pages 2022–12, 2022.

- [80] William Lotter, Gabriel Kreiman, and David Cox. Deep predictive coding networks for video prediction
 and unsupervised learning. arXiv:1605.08104, 2016.
- [81] Yanping Huang and Rajesh PN Rao. Predictive coding. Wiley Interdisciplinary Reviews: Cognitive
 Science, 2(5):580-593, 2011.
- [82] Mikail Khona and Ila R. Fiete. Attractor and integrator networks in the brain. *Nature Reviews Neuro*science, 23(12):744–766, December 2022.
- [83] Ben Sorscher, Gabriel C. Mel, Samuel A. Ocko, Lisa M. Giocomo, and Surya Ganguli. A unified theory for the computational and mechanistic origins of grid cells. *Neuron*, 111(1):121–137.e13, January 2023.
- Rylan Schaeffer, Mikail Khona, and Ila R Fiete. No free lunch from deep learning in neuroscience:
 A case study through models of the entorhinal-hippocampal circuit. In Alice H. Oh, Alekh Agarwal,
 Danielle Belgrave, and Kyunghyun Cho, editors, Advances in Neural Information Processing Systems,
 2022.
- [85] Ben Sorscher, Gabriel C. Mel, Aran Nayebi, Lisa Giocomo, Daniel Yamins, and Surya Ganguli. When and why grid cells appear or not in trained path integrators. *bioRxiv*, 2022.
- [86] 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.
- [87] D. Sussillo S. Vyas, M.D. Golub and K.V. Shenoy. Computation through neural population dynamics.
 Annual Reviews Neuroscience, 43:249–275, 2020.
- [88] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. arXiv:1301.3781 [cs. CL], 2013.
- [89] Naoki Hiratani and Haim Sompolinsky. Optimal Quadratic Binding for Relational Reasoning in Vector
 Symbolic Neural Architectures. Neural Computation, 35(2):105–155, January 2023.
- [90] Edward Paxon Frady, Denis Kleyko, and Friedrich T. Sommer. Variable binding for sparse distributed
 representations: Theory and applications. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–14, 2021.
- [91] Niru Maheswaranathan, Alex H. Williams, Matthew D. Golub, Surya Ganguli, and David Sussillo.
 Reverse engineering recurrent networks for sentiment classification reveals line attractor dynamics. Advances in neural information processing systems, 32:15696–15705, December 2019.
- [92] Mourad Mars. From word embeddings to pre-trained language models: A state-of-the-art walkthrough.

 Applied Sciences, 12(17):8805, 2022.
- [93] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
 Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learn ers. Advances in neural information processing systems, 33:1877–1901, 2020.
- [94] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online, July 2020. Association for Computational Linguistics.
- [95] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train,
 prompt, and predict: A systematic survey of prompting methods in natural language processing. ACM
 Computing Surveys, 55(9):1–35, 2023.
- [96] Micha Heilbron, Kristijan Armeni, Jan-Mathijs Schoffelen, Peter Hagoort, and Floris P De Lange. A
 hierarchy of linguistic predictions during natural language comprehension. Proceedings of the National
 Academy of Sciences, 119(32):e2201968119, 2022.

⁷⁰⁸ [97] Mante S Nieuwland. Do "early" brain responses reveal word form prediction during language comprehension? a critical review. *Neuroscience & Biobehavioral Reviews*, 96:367–400, 2019.

710 Reference Annotations

- ** [1] Investigation of whether the brain performs lazy or rich learning of task structure, in which internal representations do and do not learn task structure.
- ** [30]: Study of the dependence of the learning dynamics of a recurrent neural network on initial coupling strength. Includes a mathematical analysis of training on simple tasks predicts changes in weight structure through learning.
- ** [52]: Study of the change in representation induced by training a recurrent neural network on a classification task. Depending on initialization and task, the network expands or compresses its representation, or a combination of both; factors that encourage compression are identified empirically and through some analytical approximations of learning dynamics.
- ** [56]: Empirical study of geometric compression that occurs in the penultimate layers of image classification networks, with rich mathematical analysis of consequences. This work introduced a formal definition of compression, called neural collapse, which has inspired many follow-up studies.
- ** [11]: In-depth study of the (lack of) feature learning that occurs in the neural tangent kernel regime and beyond. This study parameterizes neural network initializations and shows how the feature learning behavior changes through a large range of these parameterizations.
- ** [69]. Study of the representations that emerge in recurrent neural networks trained to predict the future of their inputs. In this study, an agent moves randomly through an environment, and the network is trained to predict the agent's next observation. This results a spatial map of the agent's location emerging in the network representation over the course of learning.
- * [44] Study of the learning dynamics for deep linear neural networks, taking into account various possible initializations. This theory sheds light on the transition of networks from rich to lazy learning, in which internal representations do and do not learn task structure.
- * [35] Study of neural network learning that proposes that the evolution of the neural tangent kernel occurs in two phases. This assumption holds under *tabula rasa* initialization, but is shown to also be approximately true beyond the *tabula rasa* regime.