

On the Principles of Parsimony and Self-Consistency for the Emergence of Intelligence

Yi Ma^{†‡1}, Doris Tsao^{†2}, Heung-Yeung Shum^{†3}

¹Electrical Engineering and Computer Science Department, University of California, Berkeley, CA 94720, USA

²Department of Molecular & Cell Biology and Howard Hughes Medical Institute, University of California, Berkeley, CA 94720, USA

³International Digital Economy Academy, Shenzhen 518045, China

[†]E-mail: yima@eecs.berkeley.edu; dortsao@berkeley.edu; hshum@idea.edu.cn

Abstract: Ten years into the revival of deep networks and artificial intelligence, we propose a theoretical framework that sheds light on understanding deep networks within a bigger picture of Intelligence in general. We introduce two fundamental principles, *Parsimony* and *Self-consistency*, that we believe to be cornerstones for the emergence of Intelligence, artificial or natural. While these two principles have rich classical roots, we argue that they can be stated anew in entirely measurable and computable ways. More specifically, the two principles lead to an effective and efficient computational framework, compressive closed-loop transcription, that unifies and explains the evolution of modern deep networks and many artificial intelligence practices. While we mainly use modeling of visual data as an example, we believe the two principles will unify understanding of broad families of autonomous intelligent systems and provide a framework for understanding the brain.

Key words: Intelligence; Parsimony; Self-Consistency; Rate Reduction; Deep Networks; Closed-Loop Transcription

1 Context and Motivation

For an autonomous intelligent agent to survive and function in a complex environment, it must efficiently and effectively learn models that reflect both its past experiences and the current environment being perceived. Such models are critical for gathering information, making decisions, and taking actions. These models, generally referred to as world models, should be continuously improved from how projections agree with new observations and outcomes. They should incorporate both knowledge from past experiences (e.g., recognizing familiar objects) and mechanisms for interpreting immediate sensory inputs (e.g., detecting and tracking moving objects). Studies in neuroscience suggest that the brain’s world model is highly *structured* anatomically (e.g., modular brain areas and columnar organization) and functionally (e.g., sparse coding (Olshausen and Field, 1996) and subspace coding (Chang and Tsao, 2017; Bao et al., 2020)). It is believed that such a structured model is the key for the brain’s efficiency and effectiveness in perceiving, predicting, and making intelligent decisions (Barlow, 1961; Josselyn and Tonegawa, 2020).

In contrast, in the past decade, progress in artificial

intelligence has largely relied on training homogeneous black-box models, such as deep neural networks (LeCun et al., 2015), using a brute-force engineering approach. While functional modularity may emerge from training, the learned feature representation remains largely hidden or latent and difficult to interpret. As we now know, such expensive brute-force training of an end-to-end black-box model not only results in ever-growing model size and high data/computation cost¹, but is also accompanied by many problems in practice: the lack of richness in final learned representations due to neural collapse (Papayan et al., 2020)²; lack of stability in training due to mode collapse (Srivastava et al., 2017); lack of adaptiveness and susceptibility to catastrophic forgetting (McCloskey and Cohen, 1989); and lack of robustness to deformations (Azulay and Weiss, 2018; Engstrom et al., 2017) or adversarial attacks (Szegedy et al., 2013).

¹With model sizes frequently going beyond billions or trillions of parameters, even Google seems to recently have started worrying about the carbon footprint of such practices (Patterson et al., 2022)!

²This refers to the final representation for each class collapsing to a one-hot vector that carries no information about the input except its class label. Richer features might be learned inside the networks, but their structures are unclear and remain largely hidden.

[‡] Corresponding author

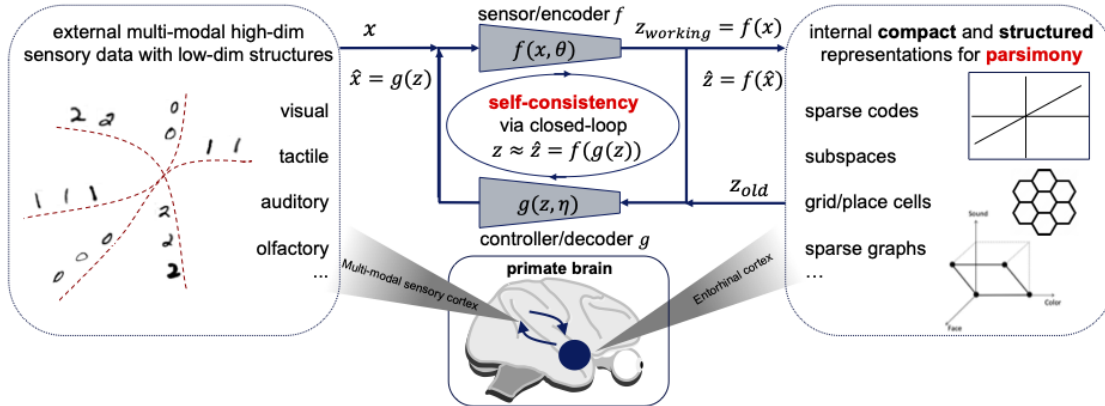


Fig. 1 Overall framework for a universal learning engine: seeking a compact and structured model for sensory data via a compressive closed-loop transcription: a (nonlinear) mapping f that maps high-dimensional sensory data with complicated low-dimensional structures to a compact structured representation. The model needs to be self-consistent in the sense that it can regenerate the original data via a map g such that f cannot distinguish despite its best effort.

A principled and unifying approach? We hypothesize that one of the fundamental reasons why these problems arise in the current practice of deep networks and artificial intelligence is a lack of systematic and integrated understanding about the functional and organizational principles of intelligent systems.

For instance, training discriminative models for classification and generative models for sampling or replaying has been largely separated in practice. Such models are typically open-loop systems that need to be trained end-to-end via supervision or self-supervision. A principle long-learned in control theory is that such open-loop systems cannot automatically correct errors in prediction, and are not adaptive to changes in the environment. This led to the introduction of “closed-loop feedback” to controlled systems so that a system can learn to correct its errors (Wiener, 1948; Mayr, 1970). As we will argue in this paper, a similar lesson can be drawn here: once a discriminative model and a generative model are combined together to form a complete closed-loop system, learning can become autonomous (without exterior supervision), and more efficient, stable and adaptive.

To understand any functional component that may be necessary in an intelligent system, such as a discriminative or a generative segment, we need to understand Intelligence from a more principled and unifying perspective. To this end, in this paper we introduce two fundamental principles: *Parsimony* and *Self-consistency*, which we believe govern the function and design of any intelligent system, artificial or natural. The two principles respectively aim to answer the following two

fundamental questions regarding learning:

1. *What to learn:* what is the objective for learning from data and how can it be measured?
2. *How to learn:* how can we achieve such an objective via efficient and effective computation?

As we will see, answers to the first question fall naturally into the realm of Information/Coding theory (Shannon, 1948) that studies how to accurately *quantify and measure* the information of the data and then to seek *the most compact* representations of the information. Once the objective of learning is clear and set, answers to the second question fall naturally into the realm of Control/Game theory (Wiener, 1948) that provides a universally effective computational framework, i.e., a *closed-loop* feedback system, for achieving any measurable objective *consistently* (Figure 1).

Basic ideas behind each of the two principles proposed in this paper can find their roots in classic works. Artificial (deep) neural networks since their earliest inception as “perceptrons” (Rosenblatt, 1958) were conceived to efficiently store and organize sensory information. Back propagation (Kelley, 1960; Rumelhart et al., 1986) was later proposed as a mechanism to learn such models. Moreover, even before the inception of neural networks, Norbert Wiener had started to contemplate computational mechanisms for learning at a system level. In his famed book *Cybernetics* (Wiener, 1948), he studied possible roles of information compression for parsimony and feedback/games in a learning machine for consistency.

But we are here to reunite and restate the two principles within the new context of data science and machine

learning, as they help better explain and unify many modern instances and practices of artificial intelligence, deep learning in particular.³ Different from earlier efforts, our restatement of these principles will be entirely *measurable* and *computationally tractable* – hence easily realizable by machines or in nature with limited resources. The purpose of this paper is to offer our overall position and perspective rather than to justify every claim technically. Nevertheless, we will provide many references to related work where readers can find convincing theoretical and compelling empirical evidence. They are based on a coherent series of past and recent developments in the study of machine learning and data science by the authors and their students (Ma et al., 2007; Wright et al., 2008; Chan et al., 2015; Yu et al., 2020; Chan et al., 2022; Baek et al., 2022; Dai et al., 2022; Tong et al., 2022; Pai et al., 2022; Wright and Ma, 2022).

Organization. In Section 2, we use visual data modeling as a concrete example to introduce the two principles and illustrate how they can be instantiated as computable objectives, architectures, and systems. In Section 3, we conjecture they lead to a universal learning engine for broader perception and decision making tasks. Finally, in Section 4, we discuss many implications of the proposed principles and their connections to neuroscience, mathematics, and higher-level intelligence.

2 Two Principles for Intelligence

In this section, we introduce and explain the two fundamental principles that can help answer the questions *what to learn* and *how to learn* by an intelligent agent or system.

2.1 What to Learn: the Principle of Parsimony

“Entities should not be multiplied unnecessarily.”

– William of Ockham

The Principle of Parsimony: *The objective of learning for an intelligent system is to identify low-dimensional structures in observations of the external world and reorganize them in the most compact and structured way.*

³As we will see, besides integrating discriminative and generative models, they lead to a closed-loop framework that works uniformly in supervised, incremental or unsupervised settings, without suffering many of the problems of open-loop deep networks.

There is a fundamental reason why intelligent systems need to embody this principle: intelligence would be impossible without it! If observations of the external world have no low-dimensional structures, there would be nothing worth learning or memorizing. There would be nothing to be relied upon for good generalization or prediction, both of which rely on new observations following the same low-dimensional structures. Thus this *principle* is not simply a *convenience* arising from the need for intelligent systems to be frugal with their resources such as energy, space, time, and matter, etc.

In some contexts, the above principle is also called the Principle of Compression. But Parsimony of Intelligence is not about achieving the best possible compression, but about *obtaining compact and structured representations via computationally efficient means*. There is no point for an intelligent system to try to compress data to the ultimate level of Kolmogorov complexity or Shannon information: they are not only intractable to compute (or even to approximate) but also result in completely unstructured representations. For example, representing data with the minimum description length (Shannon information) requires minimizing the Helmholtz free energy, say via a Helmholtz machine (Hinton et al., 1995), which is typically computationally intractable. Examined more closely, many commonly used mathematical or statistical “measures” for model goodness are either *exponentially expensive to compute* for general high-dimensional models or even become *ill-defined* for data distributions with low-dimensional supports. These measures include widely used quantities such as maximal likelihood, KL divergence, mutual information, and Jensen-Shannon and Wasserstein distances.⁴ It is commonplace in the practice of machine learning to resort to various heuristic approximations and empirical evaluations. As a result, performance guarantees and understanding are lacking.

Now we face a question: how can an intelligent system embody the Principle of Parsimony to identify and represent structures in observations in a computationally tractable and even efficient way? In theory, an intelligent system could use any family of desirable structured models for the world, as long as they are simple yet expressive enough to model useful structures in real-world sensory data. The system should be able to accurately and efficiently evaluate how good a learned model is, and the measure used should be basic, universal, and

⁴More explanations about caveats associated with these measures can be found in (Ma et al., 2007; Dai et al., 2022).

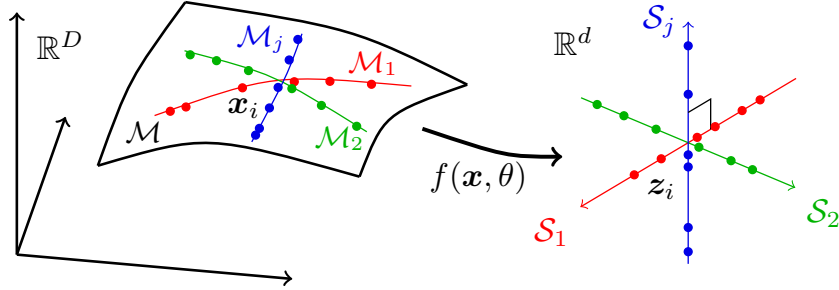


Fig. 2 Seeking a linear and discriminative representation: mapping high-dimensional sensory data, typically distributed on many nonlinear low-dimensional submanifolds, onto a set of independent linear subspaces of the same dimensions as the submanifolds.

tractable to compute and to optimize. What is a good choice for a family of structured models with such a measure?

To see how we can model and compute parsimony, we use the motivating and intuitive example of modeling visual data.⁵ To make our exposition easy, we will start with a supervised setting in this section. Nevertheless, as we will see in the next section, with parsimony as the only “self-supervision,” together with the second principle of self-consistency, a learning system can become fully autonomous and function without need for any exterior supervision.

Modeling and computing parsimony. Let us use \mathbf{x} to denote the input sensory data, say an image, and \mathbf{z} its internal representation. The sensory data sample $\mathbf{x} \in \mathbb{R}^D$ is typically rather high-dimensional (millions of pixels) but has extremely low-dimensional intrinsic structures.⁶ Without loss of generality, we may assume it is distributed on some low-dimensional submanifolds as illustrated in Figure 2. Then the purpose of learning is to establish a (usually nonlinear) mapping f , say in some parametric family $\theta \in \Theta$, from \mathbf{x} to a much lower-dimensional representation $\mathbf{z} \in \mathbb{R}^d$:

$$\mathbf{x} \in \mathbb{R}^D \xrightarrow{f(\mathbf{x}, \theta)} \mathbf{z} \in \mathbb{R}^d, \quad (1)$$

such that the distribution of features \mathbf{z} is much more compact and structured. Being compact means economic to store. Being structured implies efficient to access and use: in particular, linear structures are ideal for interpolation or extrapolation.

To be more specific and precise, we can formally instantiate the principle of Parsimony for visual data

⁵It is arguably true that among all senses, vision is the most complex to model.

⁶For example, all images of a rotating pen trace out only a one-dimensional curve in the space of millions of pixels.

modeling as trying to find a (nonlinear) transform f that achieves the following goals:

- **compression:** map high-dimensional sensory data \mathbf{x} to a low-dimensional representation \mathbf{z} ;
- **linearization:** map each class of objects distributed on a nonlinear submanifold to a linear subspace;
- **sparsification:** map different classes into subspaces with independent or maximally incoherent bases.⁷

In other words, we try to transform real-world data that may lie on a family of low-dimensional submanifolds in a high-dimensional space onto a family of independent low-dimensional linear subspaces, respectively. Such a model is called a *linear discriminative representation* (LDR) (Yu et al., 2020; Chan et al., 2022) and the compression process is illustrated in Figure 2. In some sense, one may even view the common practice of deep learning that maps each class to a “one-hot” vector as seeking a very special type of LDR models in which each target subspace is only one-dimensional and orthogonal to others.

The idea of compression as a guiding principle of the brain for representing (sensory data of) the world has strong roots in neuroscience, going back to Barlow’s efficient coding hypothesis (Barlow, 1961). Scientific studies have shown that visual object representations in the brain exhibit compact structures such as sparse codes

⁷This is related to the notion of sparse dictionary learning (Zhai et al., 2020) or independent component analysis (ICA) (Hyvärinen, 1997; Hyvärinen and Oja, 1997). Once the bases of the subspaces are made independent or incoherent by the transform, the resulting representation becomes sparse and thus collectively compact and structured. For example, two sets of subspaces with the same dimensions have the same intrinsic complexity. However, their extrinsic representations can be very different, see Figure 3. This illustrates why simply compressing data based their intrinsic complexity is not sufficient for parsimony.

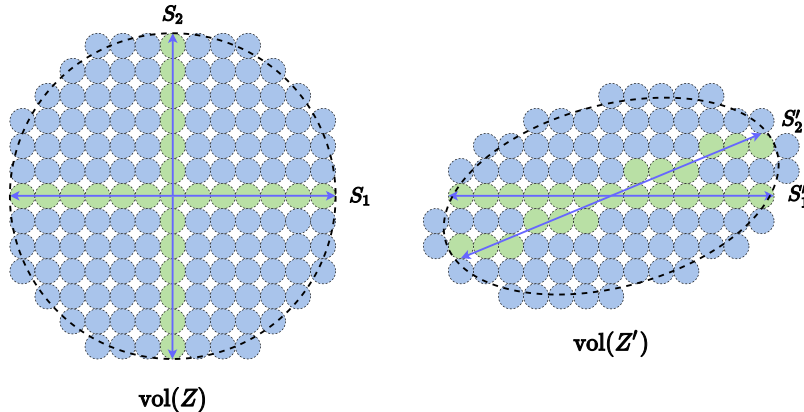


Fig. 3 Rate of all features $R = \log \#(\text{green spheres} + \text{blue spheres})$; average rate of features on the two subspaces $R^c = \log \#(\text{green spheres})$; rate reduction is the difference between the two rates: $\Delta R = R - R^c$.

(Olshausen and Field, 1996) and subspaces (Chang and Tsao, 2017; Bao et al., 2020). This supports the proposal that low-dimensional linear models are the preferred representations in the brain (at least for visual data).

Maximizing rate reduction. Somewhat remarkably, for the family of LDR models, there is a natural intrinsic measure of parsimony. Intuitively speaking, given an LDR, we can compute the total “volume” spanned by all features on all subspaces and the sum of “volumes” spanned by features of each class. Then the ratio between these two volumes gives a natural measure that suggests how good the LDR model is: the larger, the better. Figure 3 shows an example with features distributed on two subspaces S_1 and S_2 . Models on the left and right have the same intrinsic complexity. Obviously, the configuration on the left is preferred as features for different classes are made independent and orthogonal – their extrinsic representations would be the most sparse. Hence, in terms of this basic volumetric measure, the best representation should be such that “*the whole is maximally greater than the sum of its parts.*”

As per information theory, the volume of a distribution can be measured by its *rate distortion* (Cover and Thomas, 2006). Roughly speaking, the rate distortion is the logarithm of how many ϵ -balls or spheres one can pack into the space spanned by a distribution.⁸ The logarithm of the number of balls directly translates into how many binary bits one needs in order to encode a random sample drawn from the distribution subject to

⁸Sphere packing gives almost a universal way to measure the volume of space of arbitrary shape: to compare volumes of two containers, one only has to fill them both with beans and then count and compare.

the precision ϵ . This is also known more generally as the *description length* (Rissanen, 1989; Ma et al., 2007).

Now let R be the rate distortion of the joint distribution of all features $\mathbf{Z} \doteq [z^1, \dots, z^n]$ of sampled data $\mathbf{X} \doteq [x^1, \dots, x^n]$ from all, say k , classes. R^c is the average of rate distortions for the k classes: $R^c(\mathbf{Z}) = \frac{1}{k} [R(\mathbf{Z}_1) + \dots + R(\mathbf{Z}_k)]$ where $\mathbf{Z} = \mathbf{Z}_1 \cup \dots \cup \mathbf{Z}_k$. Note that, because of the logarithm, the ratio between volumes becomes the difference between rates. Then the difference between the whole and the sum of the parts, called *rate reduction* (Chan et al., 2022):

$$\Delta R(\mathbf{Z}) \doteq R(\mathbf{Z}) - R^c(\mathbf{Z}), \quad (2)$$

gives a most basic, bean-counting-like, measure for how good the feature representation \mathbf{Z} is.⁹

Although for general distributions in high-dimensional spaces the rate distortion, like many other measures mentioned before, is intractable and actually NP-hard to compute (MacDonald et al., 2019), rate distortion for data \mathbf{Z} drawn from a Gaussian supported on a subspace has a closed-form formula (Ma et al., 2007):

$$R(\mathbf{Z}) \doteq \frac{1}{2} \log \det (\mathbf{I} + \alpha \mathbf{Z} \mathbf{Z}^*). \quad (3)$$

Hence, it can be efficiently computed and optimized!

The work of Chan et al. (2022) has shown that if one uses the rate distortion functions of Gaussians and chooses a generic deep network (say a ResNet) to model the mapping $f(x, \theta)$, then by maximizing the coding rate reduction, known as *the MCR² principle*:

$$\max_{\theta} \Delta R(\mathbf{Z}(\theta)) = R(\mathbf{Z}(\theta)) - R^c(\mathbf{Z}(\theta)), \quad (4)$$

⁹The rate reduction quantity also has a natural interpretation as “information gain” (Quinlan, 1986). It measures how much information is gained, in terms of bits saved, by specifying a sample on one of the parts, compared to drawing a random sample from the whole.

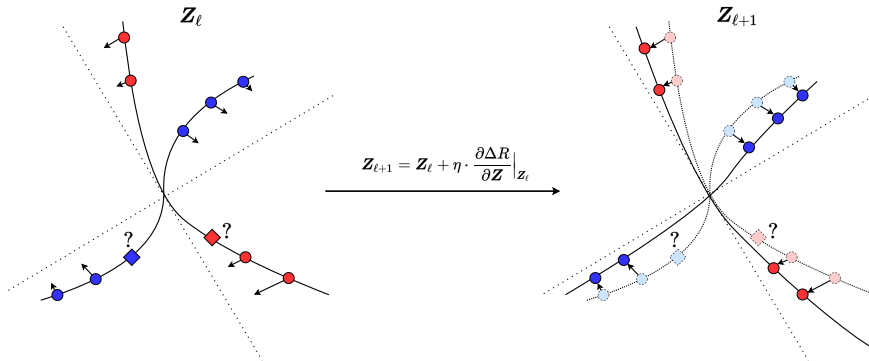


Fig. 4 A basic way to construct the nonlinear mapping f : following local gradient flow $\frac{\partial \Delta R(\mathbf{Z})}{\partial \mathbf{Z}}$ of the rate reduction ΔR , we incrementally linearize and compress features on nonlinear submanifolds and push different submanifolds apart to respective orthogonal subspaces (the two dotted lines).

one can effectively map a multi-class visual dataset to multiple orthogonal subspaces. Notice that maximizing the first term of the rate reduction R expands the volume of all features. That is, it conducts “contrastive learning” for all features simultaneously, which can be much more effective than contrasting sample pairs as normally done in conventional contrastive methods (Hadsell et al., 2006; Oord et al., 2018). Minimizing the second term R^c compresses and linearizes features in each class. This can be interpreted as conducting “contractive learning” (Rifai et al., 2011) for each class. The rate reduction objective unifies and generalizes these heuristics.

In particular, one can rigorously show that, by maximizing the rate reduction, features of different classes will be independent and features of each class will be distributed *almost uniformly* within each subspace (Chan et al., 2022). In contrast, the widely practiced *cross entropy* objective for mapping each class to a one-hot label maps final features of each class onto a one-dimensional singleton (Papayan et al., 2020).

Whitebox deep networks from unrolling optimization. Notice that in this context, the role of a deep network is simply to model the nonlinear mapping f between the external data \mathbf{x} and the internal representation \mathbf{z} . How should an intelligent system know what family of models to use for the map f in the first place? Is there a way to directly derive and construct such a mapping, instead of guessing and trying different possibilities?

Recall that our goal is to optimize the rate reduction $\Delta R(\mathbf{Z})$ as a function of the set of features \mathbf{Z} . To this end, we may directly start with the original data $\mathbf{Z}_0 = \mathbf{X}$ and optimize $\Delta R(\mathbf{Z})$ incrementally, say with a *projected*

gradient ascent (PGA) scheme:¹⁰

$$\mathbf{Z}_{\ell+1} \propto \mathbf{Z}_{\ell} + \eta \cdot \frac{\partial \Delta R}{\partial \mathbf{Z}} \Big|_{\mathbf{Z}_{\ell}} \quad \text{subject to } \|\mathbf{Z}_{\ell+1}\| = 1. \quad (5)$$

That is, one can follow the gradient of the rate reduction to move the features so that the rate reduction increases. Such a gradient-based iterative deformation process is illustrated in Figure 4.

From the closed-form formulae for the rate distortions (3), we can compute the gradient of $\Delta R = R - R^c$ in closed-form too. For example, the gradient of the first term R is of the form (similarly for the second term R^c):

$$\begin{aligned} \frac{\partial R(\mathbf{Z})}{\partial \mathbf{Z}} \Big|_{\mathbf{Z}_{\ell}} &= \frac{1}{2} \frac{\partial \log \det(\mathbf{I} + \alpha \mathbf{Z} \mathbf{Z}^*)}{\partial \mathbf{Z}} \Big|_{\mathbf{Z}_{\ell}} \quad (6) \\ &= \alpha (\mathbf{I} + \alpha \mathbf{Z}_{\ell} \mathbf{Z}_{\ell}^*)^{-1} \mathbf{Z}_{\ell} \doteq \mathbf{E}_{\ell} \mathbf{Z}_{\ell}. \quad (7) \end{aligned}$$

Similarly we can compute the gradients for the k terms $\{R(\mathbf{z}_i)\}_{i=1}^k$ in R^c and obtain k operators on \mathbf{Z}_{ℓ} , named as \mathbf{C}_i . Then the above gradient ascent operation (5) takes the following structured form:

$$\begin{aligned} \mathbf{z}_{\ell+1} &\propto \mathbf{z}_{\ell} + \eta \cdot \left[\mathbf{E}_{\ell} \mathbf{z}_{\ell} + \sigma([\mathbf{C}_{\ell}^1 \mathbf{z}_{\ell}, \dots, \mathbf{C}_{\ell}^k \mathbf{z}_{\ell}]) \right] \\ &\text{subject to } \|\mathbf{z}_{\ell+1}\| = 1, \quad (8) \end{aligned}$$

where \mathbf{E}_{ℓ} and \mathbf{C}_{ℓ} ’s are linear operators that are fully determined by covariances of the features from the previous layer \mathbf{Z}_{ℓ} (7).¹¹ Here σ is a soft-max operator that assigns \mathbf{z}_{ℓ} to its closest class based on its distance to each class, measured by $\mathbf{C}_{\ell} \mathbf{z}_{\ell}$. A diagram of all the operators per iteration is given in Figure 5 left.

¹⁰For fair comparison of coding rates between two representations, we need to normalize the scale of the features, say $\|\mathbf{z}\| = 1$.

¹¹ \mathbf{E} is associated with the gradient of the first term R and stands for “expansion” of the whole set of features, whereas \mathbf{C} ’s are associated with the gradients of multiple rate distortions in the second term R^c and stand for “compression” of features in each class. See Chan et al. (2022) for details.

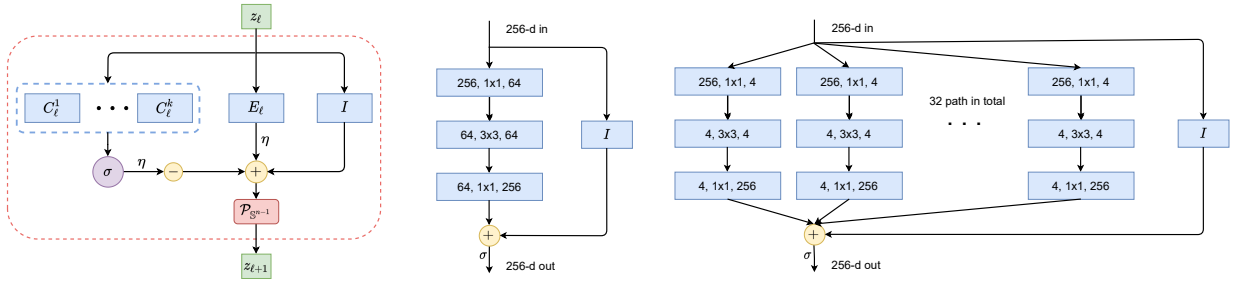


Fig. 5 Building blocks of the nonlinear mapping f . **Left:** one layer of the ReduNet as one iteration of projected gradient ascent, which precisely consists of expansive or compressive linear operators, a nonlinear softmax, plus a skip connection, and normalization. **Middle and Right:** one layer of ResNet and ResNeXt, respectively.

Acute readers may have recognized that such a diagram draws good resemblance to a layer of popular “tried-and-tested” deep networks such as ResNet (He et al., 2016) (Figure 5 middle), including parallel columns as in ResNeXt (Xie et al., 2017) (Figure 5 right) and a Mixture of Experts (MoE) (Shazeer et al., 2017). This gives a natural and principled interpretation for an important class of deep neural networks from the perspective of *unrolling an optimization scheme*. Even before the rise of modern deep networks, iterative optimization schemes for seeking sparsity such as ISTA or FISTA (Wright and Ma, 2022) had been interpreted as learnable deep networks, e.g., the work of Gregor and LeCun (2010) on Learned ISTA.¹² The class of networks derived from optimizing rate reduction has been named as *ReduNet* (Chan et al., 2022).

Forward unrolling versus backward propagation.

We see above that compression leads to an entirely constructive way of deriving a deep neural network, including its architecture and parameters, as a fully interpretable *white-box*: its layers conduct iterative and incremental optimization of a principled objective that promotes parsimony. As a result, for so-obtained deep networks, the ReduNets, starting from the data \mathbf{X} as input, each layer’s operators and parameters (E_ℓ, C_ℓ) are constructed and initialized in an entirely *forward unrolling* fashion. This is very different from the popular practice in deep learning: starting with a randomly constructed and initialized network which is then tuned globally via back propagation (Rumelhart et al., 1986). It is widely believed that the brain is unlikely to utilize back propagation as its learning mechanism due to the requirement for symmetric synapses and the complex form of the feedback. Here, the forward unrolling op-

¹²Similarly, unfolding iterative optimization for sequential sparse recovery leads to recurrent networks (Wisdom et al., 2017).

timization only relies on operations between adjacent layers that can be hard-wired; hence, it would be much easier for nature to realize and exploit.

In addition, parameters and operators of the so-constructed networks are amenable to further fine-tuning via another level of optimization, e.g., (stochastic) gradient descent realized by back propagation (Rumelhart et al., 1986).¹³ But one should not confuse the (stochastic) gradient descent used to fine-tune a network with the gradient-based optimization that layers of the network are meant to realize.

Shift-invariance and nonlinearity. If we further wish the learned encoding f to be *invariant* (or *equivariant*) to all time-shifts or space-translations, we view every sample $\mathbf{x}(t)$ with all its shifted versions $\{\mathbf{x}(t - \tau) \forall \tau\}$ as in the same equivalence class. If we compress and linearize them altogether into the same subspace, then all the linear operators E or C ’s in the above gradient operation (8) automatically become *multi-channel convolutions* (Chan et al., 2022)! As a result, the ReduNet naturally becomes a multi-channel convolution neural network (CNN), originally proposed for shift-invariant recognition (Fukushima, 1980; LeCun et al., 1998).¹⁴

Artificial selection and evolution of neural networks.

Once we realize the role of deep networks themselves is to conduct (gradient-based) iterative optimization to compress, linearize and sparsify data, it may become easy to understand the “evolution” of artificial neural

¹³It has been shown that the ReduNets have the same model capacity (say to interpolate all training data precisely) as tried-and-tested deep networks such as ResNets (Chan et al., 2022).

¹⁴In addition, due to special structures in such convolution operators E and C ’s, they are much more efficient to be computed in the *frequency domain* than in the time/space domain: the computational complexity reduces from $O(D^3)$ ¹⁵ to $O(D)$ in the dimension D of the input signals (Chan et al., 2022).

networks that has taken place in the past decade. In particular, it helps explain why only a few have emerged on top through a process of *artificial selection*: going from general MLPs to CNNs to ResNets to Transformers. In comparison, random search of network structures, such as Neural Architecture Search (Zoph and Le, 2017; Baker et al., 2017) and AutoML (Hutter et al., 2019), has not resulted in any network architectures that are effective for general tasks. We speculate that the successful architectures are simply getting more and more effective and flexible at emulating iterative optimization schemes for data compression. Besides the aforementioned similarity between ReduNet and ResNet/ResNeXt, we here discuss a few more examples.

Notice that the gradient of a rate distortion term $R(\mathbf{Z})$ is of the form (7): $\frac{\partial R}{\partial \mathbf{Z}} = \alpha(\mathbf{I} + \alpha \mathbf{Z}_\ell \mathbf{Z}_\ell^*)^{-1} \mathbf{Z}_\ell$. Instead of viewing the matrix $\alpha(\mathbf{I} + \alpha \mathbf{Z}_\ell \mathbf{Z}_\ell^*)^{-1}$ as a linear operator \mathbf{E}_ℓ acting on \mathbf{Z}_ℓ , as was done in the ReduNet, we may rewrite the whole gradient term approximately as:

$$\begin{aligned} \alpha(\mathbf{I} + \alpha \mathbf{Z}_\ell \mathbf{Z}_\ell^*)^{-1} \mathbf{Z}_\ell &\approx \alpha(\mathbf{I} - \alpha \mathbf{Z}_\ell \mathbf{Z}_\ell^*) \mathbf{Z}_\ell \\ &= \alpha[\mathbf{Z}_\ell - \alpha \mathbf{Z}_\ell (\mathbf{Z}_\ell^* \mathbf{Z}_\ell)]. \end{aligned} \quad (9)$$

That is, the gradient operation for optimizing a rate distortion term depends mostly on the auto-correlation of the features $\mathbf{A} \doteq \mathbf{Z}_\ell^* \mathbf{Z}_\ell$ from the previous iteration. This is also known as “self-attention” or “self-expression” in some contexts (Vaswani et al., 2017; Vidal, 2022). If we consider applying an additional *learnable* linear transform \mathbf{U} to each of the feature terms in the above expression (9) for the gradient, a gradient-based iteration to optimize rate distortion takes the general form:

$$\mathbf{Z}_{\ell+1} \doteq \mathbf{Z}_\ell + \mathbf{U}_o [\mathbf{Z}_\ell - \alpha \mathbf{U}_v \mathbf{Z}_\ell (\mathbf{U}_k \mathbf{Z}_\ell)^* (\mathbf{U}_q \mathbf{Z}_\ell)]. \quad (10)$$

This is of exactly the same form as the basic operation of each layer for a Transformer (Vaswani et al., 2017), also known as a self-attention (SA) head.

Moreover, very similar to ResNeXT versus ResNet, for tasks such as image classification, it is found empirically better to use not one but *multiple*, say k , SA heads in parallel in each layer (Dosovitskiy et al., 2021). In the context of rate reduction, these SA heads may be naturally interpreted as gradient terms associated with the multiple rate distortion terms in the rate reduction $\Delta R(\mathbf{Z}) = R(\mathbf{Z}) - [R(\mathbf{Z}_1) + \dots + R(\mathbf{Z}_k)]/k$. The learned linear transforms in each SA head can be interpreted as “matched filters” or “sparsifying dictionaries”¹⁶

¹⁶Interested readers may see (Zhai et al., 2020) for more details about the topic of sparse dictionary learning.

that select and transform token sets (on submanifolds) that belong to the same category (of signals or images). Hence, we conjecture that Transformers are emulating a more general family of gradient-based iterative schemes that optimize rate reduction of all the input token sets (on multiple submanifolds) by clustering, compressing, and linearizing them altogether.

Furthermore, gradient ascent or descent is merely the most basic type of optimization scheme. Networks based on unrolling such schemes (e.g., ReduNet) might not be the most efficient. One could anticipate that more advanced optimization schemes, such as accelerated gradient descent methods (Wright and Ma, 2022), could lead to more efficient deep network architectures in the future. Architecture wise, these accelerated methods require the introduction of *skip connections* across multiple layers. This may help explain, from an optimization perspective, why additional skip connections have often been found to improve network efficiency in practice, e.g., in highway networks (Srivastava et al., 2015) or dense networks (Huang et al., 2017).

2.2 How to Learn: the Principle of Self-Consistency

“Everything should be made as simple as possible, but not any simpler.”

– Albert Einstein

The principle of Parsimony alone does not ensure a learned model will capture all important information in the data sensed about the external world. For example, mapping each class to a one-dimensional “one-hot” vector, via minimizing the cross entropy, may be viewed as a form of being parsimonious. It may learn a good classifier but the features learned would collapse to a singleton, known as *neural collapse* (Papayan et al., 2020). The so learned features would no longer contain enough information to regenerate the original data. Even if we consider the more general class of LDR models, the rate reduction objective alone does not automatically determine the correct dimension of the ambient feature space. If the feature space dimension is too low, the model learned will under-fit the data; if too high, the model might over-fit.¹⁷

More generally, we take the view that perception is distinct from performance of specific tasks, and the goal of perception is to learn *everything* predictable about what is sensed. By this we mean the intelligent system should

¹⁷The first expansive or contrastive term in the rate reduction might over-expand the features to fill the space, due to noises or other variations.

be able to *regenerate the distribution of the observed data from the compressed representation* to a point that itself cannot distinguish internally despite its best effort. This view distinguishes our framework from existing frameworks that are customized to a specific class of tasks, such as the *information bottleneck* for classification (recognition) (Tishby and Zaslavsky, 2015).¹⁸ To govern the process of learning a fully faithful representation, we introduce a second principle:

The Principle of Self-consistency: *An autonomous intelligent system seeks a most self-consistent model for observations of the external world by minimizing the internal discrepancy between the observed and the regenerated.*

The two principles of Self-consistency and Parsimony are highly complementary and should always be used together. The principle of Self-consistency alone does not ensure any gain in compression or efficiency. Mathematically and computationally, it is easy and even trivial to fit any training data with over-parameterized models¹⁹ or to ensure consistency by establishing one-to-one mappings between domains with the same dimensions without learning intrinsic structures in the data distribution.²⁰ Only through compression can an intelligent system be compelled to discover intrinsic low-dimensional structures within the high-dimensional sensory data, and transform and represent them in the feature space in the most compact way for future use. Also, only through compression can we easily understand why over-parameterization, say by feature lifting with hundreds of channels as normally done in DNNs, will not lead to over-fitting if its sheer purpose is to compress in the higher-dimensional feature space: lifting helps reduce the nonlinearity in the data²¹, hence render it easier to compress and linearize.²² The role of subsequent layers

is to perform compression (and linearization), and in general the more layers, the better compressed.²³

So far we have established that we need a mechanism to determine if the compressed representation contains *all* the information that is sensed. In the remainder of this section, we will first introduce a general architecture for achieving this, a *generative model*, which can regenerate a sample from its compressed representation. A difficult problem then arises: how to sensibly measure discrepancy between the sensed sample and the regenerated sample? We argue that for an autonomous system there is one and only one solution to this, namely, measuring their discrepancy in the internal feature space. Finally, we argue that the compressive encoder and the generator must learn together through a zero-sum game. Through these deductions, we derive a universal framework for learning that we believe is inevitable.

Auto-encoding and its caveats with computability.

To ensure the learned feature mapping f and representation z have correctly captured low-dimensional structures in the data, one may check if the compressed feature z can reproduce the original data x , by some generating map g , parameterized by η :

$$x \in \mathbb{R}^D \xrightarrow{f(x,\theta)} z \in \mathbb{R}^d \xrightarrow{g(z,\eta)} \hat{x} \in \mathbb{R}^D, \quad (11)$$

in the sense that $\hat{x} = g(z, \eta)$ is close to x (according to certain measure). This process is generally known as an *auto-encoding* (Kramer, 1991; Hinton and Zemel, 1993). In the special case of compressing to a *structured* representation such as LDR, we call such an auto-encoding a *transcription*²⁴ (Dai et al., 2022). However, this goal is easier said than done. The main difficulty lies in how to make this goal computationally tractable hence physically realizable. More precisely, what is a principled measure for the difference between the distribution of x and that of \hat{x} that is both *mathematically well-defined* and *efficiently computable*? As we have mentioned before, when dealing with distributions in high-dimensional spaces with degenerate low-dimensional supports, which is almost always the case with real-world data (Ma et al., 2007; Vidal et al., 2016), conventional measures such as *must first expand*."

²³This naturally explains a seemingly mystery about deep networks: the "double-descent" phenomenon suggests a deep model's test error becomes smaller as it gets larger, after reaching its peak at certain interpolation point (Belkin et al., 2019; Yang et al., 2020).

²⁴This is analogous to the memory-forming transcription process of Engram (Josselyn and Tonegawa, 2020) or the transcription process between functional proteins and DNA (genes).

¹⁸Although in this section, for simplicity, we focus our discussions on modeling 2D imagery data, we will discuss perception of the 3D world in Section 3.1, as well as argue why perception needs to integrate recognition, reconstruction, and regeneration.

¹⁹Having a photographic memory is not intelligence. It is the same with fitting all data in the world with a Big Model.

²⁰That is the case with many popular methods for learning generative models of data such as normalizing flows (Kobyzev et al., 2021), cycle GAN (Zhu et al., 2017), and diffusion probabilistic models (Ho et al., 2020), etc. Although so learned models might be useful for applications such as image generation or style transfer, they do not identify low-dimensional structures in the data distributions nor produce compact linear structures in the learned representations.

²¹Say, as in the scattering transforms (Bruna and Mallat, 2013) or random filters (Chan et al., 2015, 2022).

²²As Lao Tzu famously said in Tao Te Ching: "That which shrinks

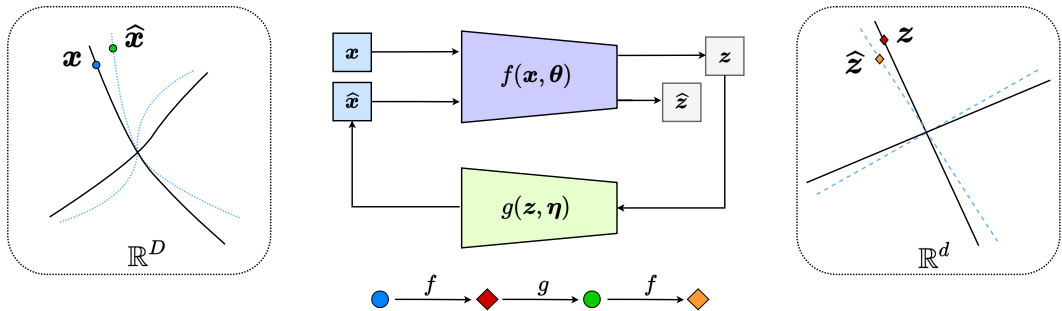


Fig. 6 A compressive closed-loop transcription of nonlinear data submanifolds to an LDR, by comparing and minimizing difference in z and \hat{z} , internally. This leads to a natural pursuit-evasion game between the encoder/sensor f and the decoder/controller g , allowing distributed of the decoded \hat{x} (the blue dotted curves) to chase and match that of the observed data x (the black solid curves).

the KL divergence, mutual information, Jensen-Shannon distance, Helmholtz free energy, and Wasserstein distances can be either ill-defined or intractable to compute, even for Gaussians (with support on subspaces) and their mixtures²⁵. How can we resolve this fundamental and yet often unacknowledged difficulty in computability associated with comparing degenerate distributions in high-dimensional spaces?

Closed-loop data transcription for self-consistency.

As we have seen in the previous section, the rate reduction ΔR gives a well-defined principled distance measure between degenerate distributions. But it is computable (with closed-form) only for mixture of subspaces or Gaussians, not for general distributions! Yet we can only expect the distribution of the internal structured representation z to be mixtures of subspaces or Gaussians, not the original data x .

This leads to a rather profound question regarding learning a “self-consistent” representation: to verify the correctness of an internal model for the external world, *does an autonomous agent really need to measure any discrepancy in the data space?* The answer is actually no. The key is to realize that, to compare x and \hat{x} , the agent only needs to compare their respective internal features $z = f(x)$ and $\hat{z} = f(\hat{x})$ via the same mapping f that intends to make z compact and structured.

$$\underline{x \xrightarrow{f(x,\theta)} z \xrightarrow{g(z,\eta)} \hat{x} \xrightarrow{f(x,\theta)} \hat{z}. \quad (12)}$$

²⁵Many existing methods formulate their objectives based on these quantities. As a result, these methods typically rely on expensive brute-force sampling to approximate these quantities or optimize their approximated lower-bounds or surrogates, such as in variational auto-encoding (VAE) (Kingma and Welling, 2013). Fundamental limitations of these methods are often disguised by good empirical results obtained with clever heuristics and excessive computational resources.

Measuring distribution difference in z space is in fact well-defined and efficient: it is arguably true that in the case of natural intelligence, learning to measure discrepancy *internally* is the only thing that the brain of a self-contained autonomous agent can do.²⁶

This effectively leads to a “closed-loop” feedback system and the overall process is illustrated as the diagram in Figure 6. The encoder f now plays an additional role as a discriminator that detects any discrepancy between x and \hat{x} through difference between their internal features z and \hat{z} . The distance between the distribution of z and that of \hat{z} can be measured through the rate reduction (2) of their samples Z and \hat{Z} :

$$\Delta R(Z(\theta), \hat{Z}(\theta, \eta)) \doteq R(Z \cup \hat{Z}) - \frac{1}{2}(R(Z) + R(\hat{Z})).$$

One can interpret popular practices for learning either a DNN classifier f or a generator g alone as learning an open-ended segment of the closed-loop system (Figure 6). This currently popular practice is very similar to an open-loop control which has long been known in the control community to be problematic and costly: training such a segment requires supervision on the desired output (e.g., class labels); and deployment of such an open-loop system is inherently not stable, robust, or adaptive if the data distributions, system parameters, or tasks change. For example, deep classification networks trained in supervised settings often suffer *catastrophic forgetting* if retrained for new tasks with new classes of data (McCloskey and Cohen, 1989). In contrast, closed-loop systems are inherently more stable and adaptive (Wiener, 1948). In fact, it has been suggested by Hinton et al. (1995) that the discriminative and generative seg-

²⁶Imagining someone colorblind, it is unlikely his/her internal representation of the world requires minimizing discrepancy in RGB values of the visual inputs x .

ments need to be combined as the “wake” and the “sleep” phases, respectively, of a complete learning process.

Self-learning through a self-critiquing game. However, just closing the loop is not enough. It is tempting to think that now we only need to optimize the generator g so as to minimize the difference between z and \hat{z} ,²⁷ say in terms of the rate reduction measure:

$$\min_{\eta} \Delta R(\mathbf{Z}(\theta), \hat{\mathbf{Z}}(\theta, \eta)). \quad (13)$$

Note that $\Delta R(\mathbf{Z}, \hat{\mathbf{Z}}) = 0$ if $\hat{\mathbf{Z}} = g(f(\mathbf{Z})) = \mathbf{Z}$. That is, the optimal learned features \mathbf{Z} should be a “fixed point” of the encoding-decoding loop²⁸. But the encoder f performs significant dimension reduction and compression, so $\hat{\mathbf{Z}} = \mathbf{Z}$ does not necessarily imply $\hat{\mathbf{X}} = \mathbf{X}$. To see this, consider the simplest case when \mathbf{X} is already on a linear subspace (say of dimension k) and f and g are a linear projection and lifting, respectively (Pai et al., 2022). f will not be able to detect any difference in its (large) null space: \mathbf{X} and any $\hat{\mathbf{X}} = \mathbf{X} + \text{null}(f)$ have the same image under f .

How can $\hat{\mathbf{Z}} = \mathbf{Z}$ imply $\hat{\mathbf{X}} = \mathbf{X}$ then? In other words, how can satisfaction of the self-consistency criterion in the internal space guarantee that we have learned to faithfully regenerate the observed data? This is possible only when the dimension k is low enough and f can be further adjusted. Let us assume the dimension of \mathbf{X} is $k < d/2$ where d is the dimension of the feature space. Then the dimension $\hat{\mathbf{X}} = g(f(\mathbf{X}))$ under a linear lifting g is a subspace of k -dimension. The union of the two subspaces of \mathbf{X} and $\hat{\mathbf{X}}$ is of dimension at most $2k < d$. Hence, if there is difference between these two subspaces and f can be an arbitrary projection, we have $f(\mathbf{X}) \neq f(\hat{\mathbf{X}})$, i.e., $\mathbf{X} \neq \hat{\mathbf{X}}$ implies $\mathbf{Z} \neq \hat{\mathbf{Z}}$.

Hence, after g minimizes the error ΔR in (13), f needs to actively adjust and detect, in its full capacity, if there is remaining discrepancy between \mathbf{X} and $\hat{\mathbf{X}}$, say by maximizing the same measure ΔR . The process can repeat between the encoder f and the decoder g and results in a natural *pursuit and evasion game*, as illustrated in Figure 6.

²⁷This is very similar in spirit to the “sleep” phase of the wake-sleep scheme proposed by Hinton et al. (1995): it essentially tries to ensure that the encoding (recognition) network f produces a response \hat{z} to the regenerated $\hat{x} = g(z)$ consistent with its origin z .

²⁸This can be viewed as a generalization to the “deep equilibrium models” (Bai et al., 2019) or the “implicit deep learning” models (El Ghaoui et al., 2021). Both interpret deep learning as conducting fixed point computation from a feedback control perspective.

In the 1961 edition of his book *Cybernetics*, Wiener (1961) had added a supplementary chapter discussing learning through playing games. The games he described were mostly about an intelligent agent against an opponent or the world (which we will discuss in the next section). Here we advocate the need of an *internal* game-like mechanism for any intelligent agent to be able to conduct self-learning via self-critique! What abides is the notion of games as a universally effective way for learning: applying the current model or strategy repeatedly against an adversarial critique, hence continuously improving the model or strategy based on feedback received through a closed loop!

Within such a framework, the encoder f assumes a dual role: in addition to learning a representation z for the data x via maximizing the rate reduction $\Delta R(\mathbf{Z})$ (as done in Section 2.1), it should also serve as a feedback “sensor” that actively detects discrepancy between the data x and the generated \hat{x} . The decoder g also assumes a dual role: it is a “controller” that corrects any discrepancy between x and \hat{x} detected by f as well as a decoder that tries to minimize the overall coding rate needed to achieve this goal (subject to a given precision).

As a result, the optimal “parsimonious” and “self-consistent” representation tuple (z, f, g) can be interpreted as the *equilibrium point* of a zero-sum game between $f(\theta)$ and $g(\eta)$, over a combined rate reduction based utility:

$$\max_{\theta} \min_{\eta} \Delta R(\mathbf{Z}(\theta)) + \Delta R(\mathbf{Z}(\theta), \hat{\mathbf{Z}}(\theta, \eta)). \quad (14)$$

Recent analysis has rigorously shown that, in the case when the input data \mathbf{X} lie on multiple linear subspaces, the desired optimal representation for \mathbf{Z} is indeed the Stackelberg equilibria (Fiez et al., 2019; Jin et al., 2019) of a sequential maximin game over a rate reduction objective similar to the above (Pai et al., 2022). It remains an open problem for the case when \mathbf{X} are on multiple non-linear submanifolds. Nevertheless, compelling empirical evidence indicates that solving this game indeed gives very good auto-encoding for real-world visual datasets (Dai et al., 2022), and automatically determines a subspace with a proper dimension for each class. It does not seem to suffer from problems such as *mode collapsing* in training conventional generative models such as GAN (Srivastava et al., 2017). The so-learned representation is simultaneously *discriminative and generative*.

Self-consistent incremental and unsupervised learning. So far we have mainly discussed the two princi-

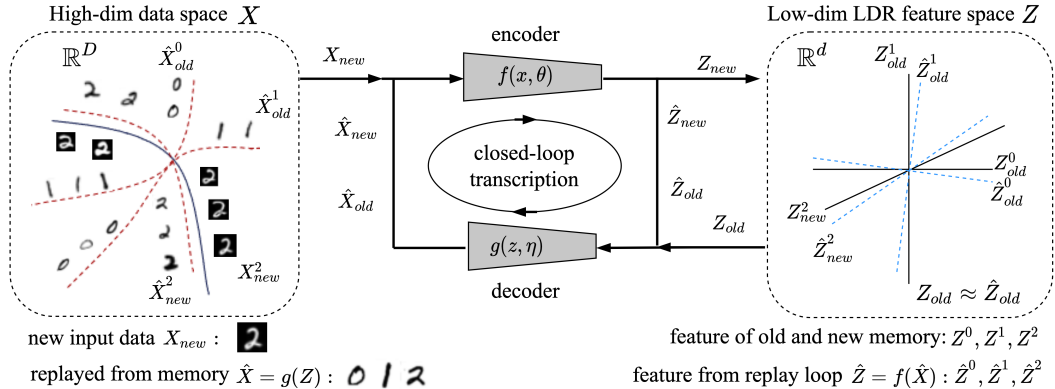


Fig. 7 Incremental learning via a compressive closed-loop transcription. For a new data class X_{new} , a new LDR memory Z_{new} is learned via a constrained minimax game between the encoder and decoder subject to a constraint that memory of past classes Z_{old} is preserved, as a “fixed point” of the closed loop.

ples in the supervised setting. In fact, one of the main advantages of our framework is that it is most natural and effective for *self-learning* via self-supervision and self-critique. In addition, since rate reduction has sought explicit (subspace-type) representations for the learned structures,²⁹ this makes it easy for past knowledge to be preserved when learning new tasks/data, as a prior (memory) to be kept *self-consistent*.

To be more clear, let us see how the closed-loop transcription framework above can be naturally extended to the case of *incremental learning* – that is, to learn to recognize one class of objects at a time instead of learning many classes jointly at the same time. While learning the representation Z_{new} for a new class, one only needs to add the cost to the objective (14) and ensure the representation Z_{old} learned before for old classes remains *self-consistent* (a fixed point) through the closed-loop transcription: $Z_{old} \approx \hat{Z}_{old} = f(g(Z_{old}))$. In other words, the above maximin game (14) becomes a constrained one:

$$\begin{aligned} \max_{\theta} \min_{\eta} \Delta R(\mathbf{Z}) + \Delta R(\mathbf{Z}, \hat{\mathbf{Z}}) + \Delta R(\mathbf{Z}_{new}, \hat{\mathbf{Z}}_{new}) \\ \text{subject to } \Delta R(\mathbf{Z}_{old}, \hat{\mathbf{Z}}_{old}) = 0. \end{aligned} \quad (15)$$

Such a constrained game makes the learning an incremental and dynamical process so that the learned transcription can continuously *adapt* to new incoming data. This process is illustrated in Figure 7.

Recent empirical studies (Tong et al., 2022) have shown that this leads to arguably the first self-contained neural system with a fixed capacity that can incrementally learn good LDR representations *without suffering*

²⁹instead of a “hidden” or “latent” representation learned by a purely generative method such as GAN (Goodfellow et al., 2014) where the features are distributed randomly in the feature space.

catastrophic forgetting (McCloskey and Cohen, 1989). The forgetting, if any at all, is rather graceful with such a closed-loop system. In addition, when images of an old class are provided again to the system to review, the learned representation can be further *consolidated* – a characteristic very similar to that of human memory. In some sense, such a constrained closed-loop formulation essentially ensures that the visual memory forming can be *Bayesian* and *adaptive* – characteristics hypothesized to be desirable for the brain (Friston, 2009).

Note that this framework is fundamentally conceived to work in an entirely unsupervised setting. Thus even though for pedagogical purposes we presented the principles assuming class information is available, the framework can be naturally extended to the entirely *unsupervised setting* in which no class information is given for any data sample. In this case, we only have to view every new sample and its augmentations as one new class in (15). This can be viewed as one type of “self-supervision.” Together with the “self-critiquing” game mechanism, a compressive closed-loop transcription can be easily learned. As shown in Figure 8, not only does the so-learned auto-encoding show good sample-wise consistency, but the learned features also demonstrate clear and meaningful local low-dimensional (thin) structures. More surprisingly, subspaces, or block-diagonal structures in feature correlation, start to emerge in features learned for the classes even without any class information provided during training at all (Figure 9)! Hence structures of so learned features resemble those of category-selective areas observed in a primate’s brain (Kanwisher et al., 1997; Kanwisher, 2010; Kriegeskorte et al., 2008; Bao et al., 2020).

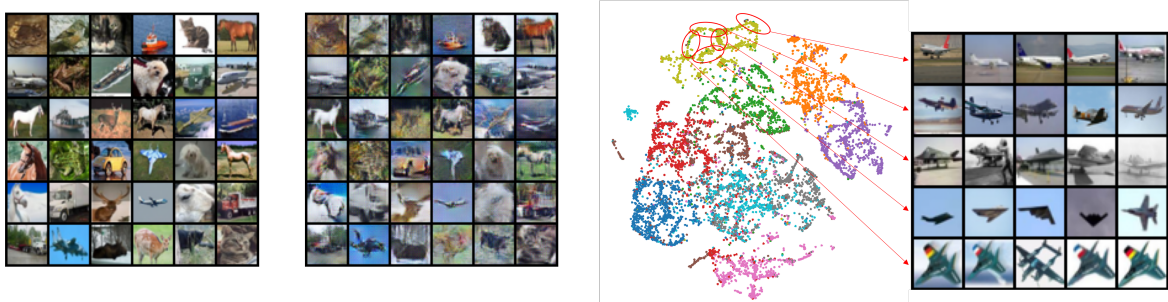


Fig. 8 Left: comparison between x and the corresponding decoded \hat{x} of the auto-encoding learned in the unsupervised setting for the CIFAR-10 dataset (with 50,000 images in 10 classes). Right: t-SNE of unsupervisedly learned features of the 10 classes and visualization of several neighborhoods with their associated images. Notice the local thin (nearly 1-D) structures in the visualized features, projected from a feature space of hundreds of dimension.

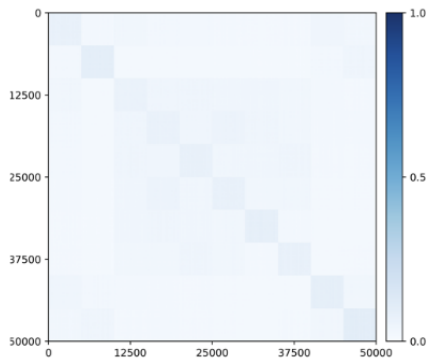


Fig. 9 Correlations between unsupervised-learned features for 50,000 images that belong to 10 classes (CIFAR-10) by the closed-loop transcription. Block-diagonal structures consistent with the classes emerge without any supervision.

3 Universal Learning Engines

“What I cannot create, I do not understand.”

– Richard Feynman

In the above section, we deduced from the first principles of Parsimony and Self-consistency the compressive closed-loop transcription framework, using the example of modeling visual imagery data. In the remaining two sections, we offer more speculative thoughts on the universality of this framework, extending it to 3D vision and reinforcement learning (the rest of this section)³⁰ and projecting its implications for neuroscience, mathematics, and higher-level intelligence (Section 4).

“Unite and build” versus “divide and conquer.”

Within the compressive closed-loop transcription framework, we have seen why and how fundamental ideas

³⁰Our discussions on the two topics require familiarity with certain domain specific terminology and knowledge. Readers who are less familiar with these topics may skip without much loss of continuity.

and concepts from coding/information theory, feedback control, deep networks, optimization, and game theory all come together to become integral parts of a complete intelligent system that can learn. Although “divide and conquer” has long been a cherished tenet in scientific research, when it comes to understanding a complex system such as Intelligence, the opposite “unite and build” should be the tenet of choice. Otherwise we would forever be *blind men with an elephant*: each person would always believe a small piece is the whole world and tend to blow its significance out of proportion.³¹

Together, the two principles serve as the glue needed to combine many necessary pieces together for the jigsaw puzzle of Intelligence, with the role of deep networks naturally and clearly revealed as models for the nonlinear mappings between external observations and internal representations. Interestingly, the principles reveal computational mechanisms for learning systems that resemble some of the key characteristics observed in or hypothesized about the brain, such as sparse coding and subspace coding (Barlow, 1961; Olshausen and Field, 1996; Chang and Tsao, 2017), closed-loop feedback (Wiener, 1948), and free energy minimization (Friston, 2009), as we will discuss more in the next section.

Notice that closed-loop compressive architectures are ubiquitous in nature for all intelligent beings and at all scales, from the brain (which compresses sensory information) to spinal circuits (which compress muscle movements) down to DNA (which compresses functional information about proteins). We believe that compressive closed-loop transcription may be the *universal learning engine* behind all intelligent behaviors. It enables intelligent beings and systems to discover and distill low-dimensional structures from seemingly complex and

³¹Hence all the superficial claims: “this or that is all you need.”

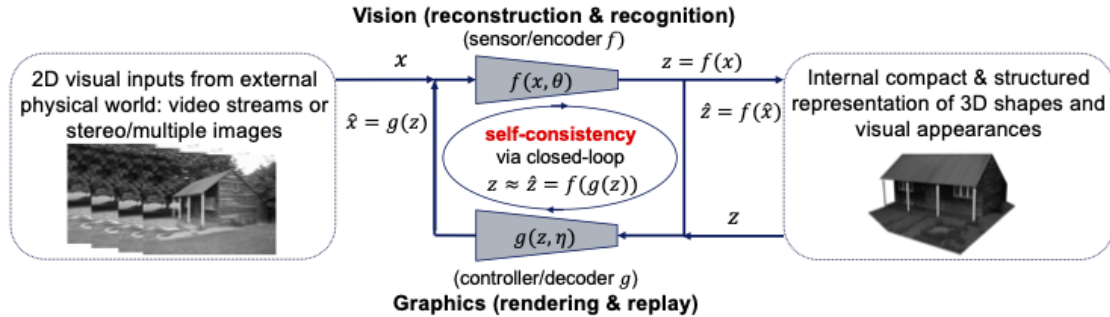


Fig. 10 A closed-loop relationship between Computer Vision and Graphics for a compact and structured 3D model of the visual inputs.

unorganized input and transform them to compact and organized internal structures to memorize and exploit.

To illustrate the universality of such a framework, for the remainder of this section, we examine two more tasks: *3D perception* and *decision making*, which are believed to be two key modules for any autonomous intelligent system (LeCun, 2022). We speculate on how, guided by the two principles, one can develop different perspectives and new insights to understand these challenging tasks.

3.1 3D Perception: Closing the Loop for Vision and Graphics

So far, we have demonstrated the success of closed-loop transcription to discover compact structures in datasets of 2D images. This has relied on the existence of *statistical correlations* among imagery data in each of the classes. We believe that the same compression mechanisms would be even more effective if the low-dimensional structures in the data are defined through hard physical or geometric constraints rather than soft statistical correlations.

In particular, if we believe the principles of Parsimony and Self-consistency also play a role in how the human brain develops mental models of the world from life-long visual inputs, then our sense of 3D space should be the result of such a closed-loop compression or transcription. The classic paradigm for 3D Vision laid out by David Marr in his influential book *Vision* (Marr, 1982) advocates a “divide and conquer” approach that partitions the task of 3D perception into several modular processes: from low-level 2D processing (e.g. edge detection, contour sketching), to mid-level 2.5D parsing (e.g. grouping, segmentation, figure and ground), and high-level 3D reconstruction (e.g. pose, shape) and recognition (e.g. objects). In contrast, the compressive closed-loop transcription proposed in this paper advocates

an opposite “unite and build” approach.

Perception as a compressive closed-loop transcription? More precisely, a three-dimensional representation of shapes, appearances, and even dynamics of objects in the world should be the most compact and structured representation that our brain has developed internally to interpret all perceived visual observations consistently. If so, the two principles then suggest that a compact and structured 3D representation is directly the internal model to be sought for. This implies that we could and should unify Computer Vision and Computer Graphics within a single closed-loop computational framework, as illustrated in Figure 10.

Computer Vision has conventionally been interpreted as a forward process that reconstructs and recognizes an internal 3D model for all the 2D visual inputs (Ma et al., 2004; Szeliski, 2022), whereas Computer Graphics (Hughes et al., 2014) represents its inverse process that renders and animates the internal 3D model. There might be tremendous computational and practical benefits to combine these two processes directly into a closed-loop system: all the rich structures (e.g. sparsity and smoothness) in geometric shapes, visual appearances, and dynamics can be exploited together for a unified 3D model that is the most compact and consistent with all the visual inputs.

Indeed, the recognition techniques in computer vision could help computer graphics in building compact models in the spaces of shapes and appearance and enabling new ways for creating realistic 3D content. On the other hand, the 3D modeling and simulation techniques in computer graphics could predict, learn and verify the properties and behaviors of the real objects and scenes analyzed by computer vision algorithms. In fact, the approach of “analysis by synthesis” has been long practiced by the vision and graphics community, e.g., for efficient

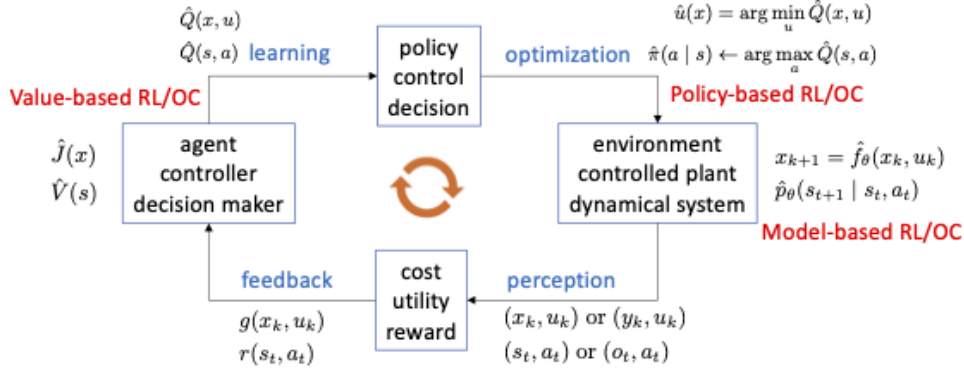


Fig. 11 An autonomous intelligent agent that integrates perception (feedback), learning, optimization, and action in a closed loop to learn optimal policy for a certain task. s_t or x_k is the state of the world model; r or g is the perceived reward or cost of action a_t or control u_k on the current state; J or V is the (learned) cost or value associated with each state, Q is the (learned) cost associated with each state-action pair. Here, we deliberately use terminologies from optimal control (OC) (Bertsekas, 2012) and reinforcement learning (RL) (Sutton and Barto, 2018) in parallel, for both comparison and unification.

online perception (Yildirim et al., 2020). Some recent examples of closing the loop for computer vision and graphics include a learned 3D rendering engine (Kulkarni et al., 2015) and 3D aware image synthesis (Chan et al., 2020; Wood et al., 2021).

Unified representations for appearance and shape?

Image-based rendering (Levoy and Hanrahan, 1996; Gortler et al., 1996; Shum et al., 2011), in which a new view is generated by learning from a set of given images, may be regarded as an early attempt to close the gap between vision and graphics with the principles of parsimony and self-consistency. In particular, plenoptic sampling (Chai et al., 2000) showed that an anti-aliased image (self-consistency) can be achieved with the minimum number of images required (parsimony).

Recent developments in modeling radiance fields have provided more empirical evidence for this view (Yu et al., 2021): directly exploiting low-dimensional structures in the radiance field in 3D (sparse support and spatial smoothness) leads to much more efficient and effective solutions than brute-force training of black-box deep neural networks (Mildenhall et al., 2020). However, it remains a challenge for the future to identify the right family of compact and structured 3D representations that can integrate shape geometry, appearance, and even dynamics in a unified framework that leads to the minimal complexity in data, model, and computation.

3.2 Decision Making: Closing the Loop for Perception, Learning, and Action

So far in this paper, we have discussed how compressive closed-loop transcription may lead to an effective

and efficient framework for learning a good perceptual model from visual inputs. At the next level, such a perceptual model can then be used by an autonomous agent to achieve certain tasks in a complex *dynamical* environment. The overall process for the agent to learn from perceived results or received rewards of its actions forms another closed loop at a higher level (Figure 11).

The principle of self-consistency is clearly at play here: the role of the closed-loop feedback system is to ensure that the learned model and control policy by the agent is *consistent* with the external world in such a way that the model can make the best prediction of the state (s_t) transition, and the learned control policy π_θ for the action (a_t) results in maximal expected reward R :

$$\max_{\theta} R(\theta) \doteq \mathbb{E}_{a_t \sim \pi_\theta(s_t)} \left[\sum_t r(s_t, a_t) \right]. \quad (16)$$

Note here the reward R plays a similar role as the rate reduction objective (4) for LDR models, which measures the “goodness” of the learned control policy π and guides its improvement.

The principle of Parsimony is the main reason for the very success of modern reinforcement learning in tackling large-scale tasks such as Alpha-Go (Silver et al., 2016, 2017) and playing video games (Berner et al., 2019; Vinyals et al., 2019). In almost all tasks that have a state-action space of astronomical size or dimension, say D , practitioners always assume the optimal value function V^* , Q-function Q^* , or policy π^* only depends on a small number of, say $d \ll D$, features:

$$\begin{aligned} V^*(s) &\approx \hat{V}(f(s, a)), \\ Q^*(s, a) &\approx \hat{Q}(f(s, a)), \\ \pi^*(a | s) &\approx \hat{\pi}(a; f(s, a)), \end{aligned} \quad (17)$$

where $f(\mathbf{s}, \mathbf{a}) \in \mathbb{R}^d$ is a nonlinear mapping that learns some low-dimensional features of the extremely large or high-dimensional state-action space. In the case of video games, the state dimension D is easily in the millions and yet the number of features d needed to learn a good policy is typically only a few dozen or hundred! Very often, these optimal control policy or value/reward functions sought in OC/RL are even assumed to be a linear superposition of these features (Ng and Russell, 2000; Kakade, 2001):

$$\omega^\top f(\mathbf{s}, \mathbf{a}) = \omega_1 \cdot f_1(\mathbf{s}, \mathbf{a}) + \dots + \omega_d \cdot f_d(\mathbf{s}, \mathbf{a}). \quad (18)$$

That is, the nonlinear mapping f is assumed to be able to also *linearize* the dependency of the policy/value/reward functions on the learned features.³²

Autonomous feature selection via a game? Notice that all these practices in RL are very similar in spirit to the learning objectives under the principle of Parsimony stated in Section 2.1. Effectively exploiting the low-dimensional structures is the (only) reason why the learning can be so *scalable* with such a high-dimensional state-action space; and correctly identifying and linearizing such low-dimensional structures is the key for the so-learned control policy to be *generalizable*³³. Nevertheless, a proper choice in the number of features d remains heuristically designed by human in practice. That makes the overall RL not autonomous. We believe that, for a closed-loop learning system to automatically determine the right number of features associated with a reward/task, one has to extend the formulation of RL (16) to a certain maximin game³⁴, in a similar spirit as those studied in Section 2.2 for achieving Self-consistency for visual modeling.

Data and computational efficiency of RL? In recent years, there have been many theoretical attempts to explain the empirically observed efficiency of reinforcement learning in terms of sampling and computation complexity of Markov decision processes (MDP). However, any theory based on unstructured generic MDPs and reward functions would not be able to provide pertinent explanations to such empirical successes. For example, some of

³²It has been a common practice in systems theory to linearize any nonlinear dynamics before controlling them, through either nonlinear mappings known as the Koopman operators (Koopman, 1931) or feedback linearization (Sastry, 1999).

³³otherwise, the learned model/policy tends to over-fit or under-fit.

³⁴by introducing a self-critique for the features selected and learned.

the best known bounds on the sample complexity for reinforcement learning remains linear in cardinality of the state space and action $O(|\mathcal{S}||\mathcal{A}|)$ (Li et al., 2020), which does not explain the empirically observed efficiency of RL in large-scale tasks (such as Alpha-Go and video games) where the state or action spaces are astronomical.

We believe that the efficiency of RL in tackling many of the practical large-scale tasks can come *only* from the intrinsic low-dimensionality in the system dynamics or correlation between the optimal policy/control and the states. For example, assuming the systems have bounded eluder dimension (Osband and Roy, 2014) or the MDPs are low rank (Uehara et al., 2021; Agarwal et al., 2020). The role of deep networks is again to identify and model such low-dimensional structure and hopefully linearize it.

To conclude, for large-scale RL tasks, it is the two principles together that make such a closed-loop system of perception, learning, and action a truly efficient and effective learning engine. With such an engine, autonomous agents are able to discover low-dimensional structures if there are indeed such structures in the environment and in the learning task, and eventually act intelligently when the structures learned are good enough and generalize well!

4 A Broader Program for Intelligence

“If I were to choose a patron saint for cybernetics out of the history of science, I should have to choose Leibniz. The philosophy of Leibniz centers about two closely related concepts – that of a universal symbolism and that of a calculus of reasoning.”

– Norbert Wiener, *Cybernetics*, 1961

It has been ten years since the dramatic revival of deep neural networks with the work of Krizhevsky et al. (2012), which has led to considerable enthusiasm about artificial intelligence in both the technology industry and the scientific community. Subsequent theoretical studies of deep learning often view deep networks themselves as the object of study (Roberts et al., 2022). We here however have argued that deep networks are better understood as a means to an end: they clearly arise to serve the purposes of identifying and transforming nonlinear low-dimensional structures in high-dimensional data, a universal task for learning from high-dimensional data (Wright and Ma, 2022).

More broadly, in this paper, we have proposed and

argued that Parsimony and Self-consistency are two fundamental principles responsible for the emergence of Intelligence, artificial or natural. The two principles together lead to a closed-loop computational framework that unifies and clarifies many practices and empirical findings of deep learning and artificial intelligence. Furthermore, we believe they will guide us from now on to study Intelligence with a more principled and integrated approach. Only thus can we achieve a new level in understanding the science and mathematics of Intelligence.

4.1 Neuroscience of Intelligence

One would naturally expect any fundamental principle of intelligence to have major implications for the design of the most intelligent thing in the universe, the brain. As already mentioned, the principles of Parsimony and Self-consistency shed new light on several experimental observations concerning the primate visual system. Even more importantly, they shine light on what to look for in future experiments.

We have shown that seeking an internally parsimonious and predictive representation alone is enough “self-supervision” to allow structures to emerge automatically in the final representation learned through a compressive closed-loop transcription. For example, Figure 9 shows that unsupervised data transcription learns features that automatically distinguish different categories, providing an explanation for category-selective representations observed in the brain (Kanwisher et al., 1997; Kanwisher, 2010; Kriegeskorte et al., 2008; Bao et al., 2020). These features also provide a plausible explanation for the widespread observations of sparse coding (Olshausen and Field, 1996) and subspace coding (Chang and Tsao, 2017; Bao et al., 2020) in the primate brain. Furthermore, beyond modeling of visual data, recent studies in neuroscience suggest the emergence of other structured representations in the brain such as “place cells” might also be the result of coding spatial information in the most compressed way (Benna and Fusi, 2021).

Arguably, the maximal coding rate reduction (MCR²) principle (4) is similar in spirit to the “free energy minimization principle” from cognitive science (Friston, 2009), which attempts to provide a framework for Bayesian inference through energy minimization. But unlike the general notion of free energy, rate reduction is computationally tractable and directly optimizable as it can be expressed in closed form. Furthermore, the interplay of our two principles suggests that autonomous learning of the correct model (class) should be done

via a closed-loop maximin game over such a utility (14), instead of minimization alone. Thus we believe our framework provides a new perspective on how to practically implement Bayesian inference.

Our framework clarifies the overall learning architecture used by the brain. One important insight is that a feedforward segment can be constructed via unrolling an optimization scheme; learning from a random network via back propagation is not necessary. Furthermore, our framework suggests the existence of a complementary generative segment to form a closed-loop feedback system to guide learning. Finally, our framework sheds light on the elusive “prediction error” signal sought by many neuroscientists interested in brain mechanisms for “predictive coding,” a computational scheme with resonances to compressive closed-loop transcription (Rao and Ballard, 1999; Keller and Mrcic-Flogel, 2018): for reasons of computational tractability, the discrepancy between incoming and generated observations should be measured at the final stage of representation.

So far, many of the resemblances between this new framework and natural intelligence discussed in this paper are still speculative and remain to be corroborated with future scientific experiments and evidences. Nevertheless, these speculations offer neuroscientists ample new perspectives and hypotheses about why and how Intelligence could emerge in nature.

4.2 Mathematics of Intelligence

In terms of mathematical or statistical models for data analysis, one can view our framework as a generalization of PCA (Jolliffe, 1986), GPCA (Vidal et al., 2016), RPCA (Candès et al., 2011), and Nonlinear PCA (Kramer, 1991) to the case with multiple low-dimensional nonlinear submanifolds in a high-dimensional space. These classical methods largely model data with a single or multiple linear subspaces or with a single nonlinear submanifold. We have argued that the role of deep networks is mainly to model the nonlinear mappings that linearize and separate multiple low-dimensional submanifolds simultaneously. This generalization is necessary to connect these idealistic classic models to true structures of real-world data. Despite promising and exciting empirical evidence, mathematical properties of the compressive closed-loop transcription process have not been so well studied and understood. To our best knowledge, only for the case when the original data x are assumed to be on multiple linear subspaces, it has been shown the maximin game based on rate reduction gives the correct

optimal solution (Pai et al., 2022). Little is known about transcription of nonlinear submanifolds.

A rigorous and systematic investigation requires understanding of high-dimensional probability distributions with low-dimensional supports on submanifolds (Fefferman et al., 2013). Hence mathematically, it is crucial to study how such submanifolds in high-dimensional spaces can be identified, grouped or separated, deformed and flattened with minimal distortion to their original metric, geometry, and topology (Tenenbaum et al., 2000; Buchanan et al., 2020; Wang et al., 2021). Problems like these seem to fall in an understudied area between classical Differential Geometry and Differential Topology in mathematics.

In addition, often we also wish that during the process of deformation, the probability measure of data on each submanifold is redistributed in a certain optimal way such that coding and sampling will be the most economic and efficient. This is related to topics such as Optimal Transport (Lei et al., 2017). For the case when the submanifolds are fixed linear subspaces, understanding the achievable extremals of the rate reduction, or ratio of volumes of the whole and the parts, seems related to certain fundamental inequalities in analysis such as the *Brascamp-Lieb* inequalities (Jonathan Bennett et al., 2007). More general problems like these seem to be related to the studies of Metric Entropy and Coding Theory for distributions over more general compact structures.

Besides nonlinear low-dimensional structures, real-world data and signals are typically *invariant* to shift in time, translation in space, or to more general group transformations. Wiener (1961) recognized dealing both nonlinearity and invariance simultaneously presents a major technical challenge. He had made early attempts to generalize Harmonic Analysis to nonlinear processes and systems.³⁵ Indeed, the revival of deep learning has reignited strong interests in this important problem and great progresses have been made recently, including the work of Bruna and Mallat (2013); Wiatowski and Bölcskei (2018); Cohen and Welling (2016); Cohen et al. (2019). Our framework suggests that a more unifying approach to deal with nonlinearity and invariance is through (incremental) compression. That has led to a natural derivation of structured deep networks such as the convolution ReduNet (Chan et al., 2022). We believe compression provides a unifying perspective to modeling general nonlinear sequential data and processes too, which could lead to mathematically rigorous justification

³⁵He used his analysis to explain brain waves (Wiener, 1961)!

for popular models such as RNNs or LSTMs (Hochreiter and Schmidhuber, 1997).

But besides pure mathematical interests, we must require this time that the mathematical investigation leads to *computationally tractable* measures and *scalable* algorithms. One has to characterize the precise statistical and computational resources needed to achieve such tasks, in the same spirit as the research program set for Compressive Sensing (Wright and Ma, 2022), because Intelligence needs to apply them to model high-dimensional data and solve large-scale tasks. This entails to “close the loop” between Mathematics and Computation, enabling the use of rich families of good geometric structures (e.g., sparse models, subspaces, grids, groups, or graphs; Figure 1, right) as compact archetypes for modeling real-world data, through efficiently computable nonlinear mappings that generalize deep networks, see e.g. Bronstein et al. (2021).

4.3 Towards Higher-Level Intelligence

The two principles laid out in this paper are mainly for explaining the emergence of Intelligence in *individual* agents or systems, related to the notion of *ontogenetic learning* that Nobert Wiener first proposed (Wiener, 1948). It is probably not a coincidence that after more than seventy years, we find ourselves in this paper “closing the loop” of modern practice of Artificial Intelligence back to its roots in Cybernetics and interweaving the very same set of fundamental concepts that Wiener touched upon in his book while conducting inquiries into the jigsaw puzzle of Intelligence: *information and compression, closed-loop feedback, learning via games, white-box modeling, nonlinearity, shift-invariance, etc.*

As we see from this paper, the compressive closed-loop transcription is arguably the first computational framework that integrates many of these pieces coherently together. It is in close spirit to earlier frameworks (Hinton et al., 1995) but makes them computationally tractable and scalable. In particular, the learned nonlinear encoding/decoding mappings of the loop, often manifested as deep networks, essentially provide a much needed “interface” between external unorganized raw sensory data (say visual, auditory etc.) and internal compact and structured representations.

However, the two principles proposed in this paper do not necessarily explain all aspects of Intelligence. Obviously, computational mechanisms behind the emergence and development of high-level semantic, symbolic, or logic inferences remain elusive, although many foun-

dational works have been set forth by pioneers like John McCarthy, Marvin Minsky, Allen Newell and Herbert Simon since the 1950's (Simon, 1969; Newell and Simon, 1972) and a comprehensive modern exposition can be found in Russell and Norvig (2020). Till today, there are still active and contentious debates about whether such high-level symbolic intelligence can emerge from continuous learning or has to be hard-coded (Marcus, 2020; LeCun and Browning, 2022).

In our view, structured internal representations such as subspaces are one necessary intermediate step for the emergence of high-level semantic or symbolic concepts—each subspace corresponds to one *discrete* category (of objects). Additional statistical, causal or logical relationships among the so-abstracted discrete concepts can be further modeled parsimoniously as a compact and structured (say sparse) graph, with each node representing a subspace/category. The graph can be learned via auto-encoding to ensure self-consistency, e.g. (Bear et al., 2020).

We conjecture that only on top of *compact and structured* representations learned by individual agents can the emergence and development of high-level intelligence (with shareable symbolic knowledge) be possible, subsequently and eventually. We suggest that new principles for the emergence of high-level intelligence, if any, should be sought through the need for efficient communication of information or transfer of knowledge among intelligent agents. This is related to the notion of *phylogenetic learning* that Wiener also discussed (Wiener, 1961).

Furthermore, any new principles needed for such higher-level intelligence must reveal *reasons* why alignment and sharing of internal models/concepts across different individual agents is possible computationally, as well as reveal certain measurable *gains* in intelligence for a group of agents from such symbolic abstraction and sharing.

Intelligence as interpretable and computable systems.

Obviously, as we advance our inquiries into higher-level Intelligence, we want to set much higher standards this time. Whatever new principles there might remain to be discovered in the future, for them to truly play a substantial role in the emergence and development of Intelligence, they must share two characteristics to the two principles we have presented in this paper:

1. **Interpretability:** *all principles together should help reveal computational mechanisms of Intelli-*

*gence as a white box*³⁶, including measurable objectives, associated computing architectures, and structures of learned representations.

2. **Computability:** *any new principle for Intelligence must be computationally tractable and scalable, physically realizable by computing machines or nature, and ultimately corroborated with scientific evidence.*

Only with such fully interpretable and truly realizable principles in place can we explain all existing intelligent systems, artificial or natural, as partial or holistic instantiations of these principles. Then they can help us to discover effective architectures and systems for different intelligent tasks without relying on the current expensive and time-consuming “trial-and-error” approach to advance. We will also be able to characterize the minimal data and computation resources needed to achieve these tasks, instead of the current brute-force approach that simply advocates “the bigger, the better.” Intelligence should not be the privilege of the most resourceful, as it is not the way of nature. Instead, parsimony and autonomy are the main characteristics.³⁷ Under a correct set of principles, anyone should be able to design and build future generations of intelligent systems, small or big, with autonomy, ability and efficiency emulating and even surpassing that of animals and humans, eventually.

5 Conclusion

Through this paper, we hope to have convinced the reader that we are now at a much better place than people like Wiener and Shannon seventy years ago when it comes to uncovering, understanding, and even exploiting the works of Intelligence. We have proposed and argued that, under the two principles of Parsimony and Self-consistency, it is possible to assemble many necessary pieces of the puzzle of Intelligence into a unified computational framework that is easily implementable on machines or by nature. This unifying framework offers new perspectives on how we could further advance the study of perception, decision making, and intelligence in general.

To draw an end to our proposal for a principled approach to Intelligence, we emphasize once again that all

³⁶Again, the phrase “white box” modeling has been conveniently borrowed from Wiener’s *Cybernetics* (Wiener, 1961).

³⁷A tiny ant is arguably much more intelligent and independent than any legged robot in the world, which has merely a quarter of million neurons consuming negligible energy.

scientific principles for Intelligence should not be merely philosophical guidelines or merely conceptual frameworks formulated with mathematical/statistical quantities that are intractable to compute or can only be approximated heuristically. They should rely on the most basic and principled objectives that are measurable with finite observations and lead to computational systems that can be realized even with limited resources. This belief is probably best expressed through a quote from Lord Kelvin:³⁸

“When you can measure what you are speaking about and express it in numbers, you know something about it; but when you cannot measure it, when you cannot express it in numbers, your knowledge is of the meager and unsatisfactory kind: it may be the beginning of knowledge, but you have scarcely, in your thoughts, advanced to the stage of science, whatever the matter may be.”

– Lord Kelvin, 1883

Afterword and Acknowledgement

Although the research of Yi and Harry focuses more on computer vision and computer graphics, they both happened to major in control and automation as undergraduate students. They started their collaboration many years ago at Microsoft Research Asia (MSRA) with a compression-based approach to data clustering and classification (Ma et al., 2007; Wright et al., 2008). In the past two years, they have had frequent discussions and debates about understanding (deep) learning and (artificial) intelligence. Their shared interests in Intelligence have brought all these fundamental ideas together and led to the recent collaboration on closed-loop transcription (Dai et al., 2022), and eventually to many of the views shared in this paper. Doris is deeply interested in whether and how the brain implements generative models for visual perception, and her group has been having intense discussions with Yi on this topic since moving to UC Berkeley a year ago.

The idea of writing of this position paper is partly motivated from recent stimulating discussions among a group of researchers with very diverse backgrounds in artificial intelligence, applied mathematics, optimization, and neuroscience: Professor John Wright and Stefano

Fusi of Columbia University, Professor Yann LeCun and Rob Fergus of New York University, Dr. Xin Tong of MSRA. We realize that these perspectives might be interesting to broader scientific and engineering communities.

Some of the thoughts presented about integrating pieces of the puzzle for intelligent systems can be traced back to an advanced robotics course that Yi had led and organized jointly with Professor Jitendra Malik, Shankar Sastry, and Claire Tomlin as Berkeley EECS290-005: *the Integration of Perception, Learning, and Control* in Spring 2021. The need for an integrated view or a “unite and build” approach seems to be a topic that is drawing increasing interest and importance for the study of Artificial Intelligence.

We would also like to thank many of our former and current students who, against extraordinary odds, have worked on projects under this new framework in the past several years when some of the ideas were still at their infancy and seemed not in accordance with the mainstream, including Xili Dai, Yaodong Yu, Peter Tong, Ryan Chan, Chong You, Ziyang Wu, Christina Baek, Druv Pai, Brent Yi, Michael Psenska, and others. Many of the technical evidences and figures used in this position paper are conveniently borrowed from their recent research results.

References

- Agarwal A, Kakade S, Krishnamurthy A, et al., 2020. Flambe: Structural complexity and representation learning of low rank MDPs. *Advances in neural information processing systems*, 33:20095-20107.
- Azulay A, Weiss Y, 2018. Why do deep convolutional networks generalize so poorly to small image transformations? *arXiv preprint arXiv:180512177*.
- Baek C, Wu Z, Ryan Chan TD, et al., 2022. Efficient maximal coding rate reduction by variational form. *CVPR*.
- Bai S, Kolter JZ, Koltun V, 2019. Deep equilibrium models. *Advances in Neural Information Processing Systems*, 32.
- Baker B, Gupta O, Naik N, et al., 2017. Designing neural network architectures using reinforcement learning. *ArXiv*, abs/1611.02167.
- Bao P, She L, McGill M, et al., 2020. A map of object space in primate inferotemporal cortex. *Nature*, 583:103-108.
- Barlow HB, 1961. Possible principles underlying the transformation of sensory messages. *Sensory communication*, 1(01).
- Bear D, Fan C, Mrowca D, et al., 2020. Learning physical graph representations from visual scenes. *Advances in Neural Information Processing Systems*, 33.
- Belkin M, Hsu D, Ma S, et al., 2019. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849-15854.
- Benna MK, Fusi S, 2021. Place cells may simply be memory cells: Memory compression leads to spatial tuning and history dependence. *Proceedings of the National Academy of Sciences*, 118(51).

³⁸that we have learned from Professor Jitendra Malik of UC Berkeley.

- Berner C, Brockman G, Chan B, et al., 2019. Dota 2 with large scale deep reinforcement learning. *preprint arXiv:191206680*, .
- Bertsekas D, 2012. Dynamic Programming and Optimal Control, volume I and II. Athena Scientific.
- Bronstein MM, Bruna J, Cohen T, et al., 2021. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges, . <http://arxiv.org/abs/2104.13478>
- Bruna J, Mallat S, 2013. Invariant scattering convolution networks. *PAMI*, .
- Buchanan S, Gilboa D, Wright J, 2020. Deep networks and the multiple manifold problem. *International Conference on Learning Representations*, . <https://par.nsf.gov/biblio/10218695>
- Candès EJ, Li X, Ma Y, et al., 2011. Robust principal component analysis? *Journal of the ACM (JACM)*, 58(3):1-37.
- Chai JX, Tong X, Chan SC, et al., 2000. Plenoptic sampling. Proceedings of the 27th annual conference on Computer graphics and interactive techniques, p.307-318.
- Chan E, Monteiro M, Kellnhofer P, et al., 2020. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. *arXiv*.
- Chan KHR, Yu Y, You C, et al., 2022. ReduNet: A white-box deep network from the principle of maximizing rate reduction. *Journal of Machine Learning Research*, 23(114):1-103. <http://jmlr.org/papers/v23/21-0631.html>
- Chan TH, Jia K, Gao S, et al., 2015. PCANet: A simple deep learning baseline for image classification? *IEEE Transactions on Image Processing*, .
- Chang L, Tsao D, 2017. The code for facial identity in the primate brain. *Cell*, .
- Cohen T, Geiger M, Weiler M, 2019. A general theory of equivariant CNNs on homogeneous spaces. *NeurIPS*.
- Cohen TS, Welling M, 2016. Group equivariant convolutional networks. *CoRR*, abs/1602.07576. <http://arxiv.org/abs/1602.07576>
- Cover TM, Thomas JA, 2006. Elements of Information Theory. Wiley-Interscience.
- Dai X, Tong S, Li M, et al., 2022. CTRL: Closed-loop data transcription to an LDR via minimaxing rate reduction. *Entropy*, *arXiv:211106636*, .
- Dosovitskiy A, Beyer L, Kolesnikov A, et al., 2021. An image is worth 16x16 words: Transformers for image recognition at scale. International Conference on Learning Representations, <https://openreview.net/forum?id=YicbFdNTTy>
- El Ghaoui L, Gu F, Travacca B, et al., 2021. Implicit deep learning. *SIAM Journal on Mathematics of Data Science*, 3(3):930-958. <https://doi.org/10.1137/20M1358517> <https://doi.org/10.1137/20M1358517>
- Engstrom L, Tran B, Tsipras D, et al., 2017. A rotation and a translation suffice: Fooling CNNs with simple transformations. *arXiv preprint arXiv:171202779*, .
- Fefferman C, Mitter S, Narayanan H, 2013. Testing the manifold hypothesis. *preprint arXiv:13100425*, . <https://arxiv.org/abs/1310.0425>
- Fiez T, Chasnov B, Ratliff LJ, 2019. Convergence of Learning Dynamics in Stackelberg Games. *preprint arXiv:190601217*, . <https://arxiv.org/abs/1906.01217>
- Friston K, 2009. The free-energy principle: a rough guide to the brain? *Trends in Cognitive Sciences*, 13(7):293-301. <https://doi.org/https://doi.org/10.1016/j.tics.2009.04.005>
- Fukushima K, 1980. Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biol Cybernetics*, 36:193-202.
- Goodfellow I, Pouget-Abadie J, Mirza M, et al., 2014. Generative adversarial nets. *Advances in neural information processing systems*, 27.
- Gortler SJ, Grzeszczuk R, Szeliski R, et al., 1996. The lumigraph. Proceedings of the 23rd annual conference on Computer graphics and interactive techniques, p.43-54.
- Gregor K, LeCun Y, 2010. Learning fast approximations of sparse coding. Proceedings of the 27th International Conference on International Conference on Machine Learning, p.399-406.
- Hadsell R, Chopra S, LeCun Y, 2006. Dimensionality reduction by learning an invariant mapping. 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, CVPR 2006, p.1735-1742.
- He K, Zhang X, Ren S, et al., 2016. Deep residual learning for image recognition. Proceedings of the IEEE conference on computer vision and pattern recognition, p.770-778.
- Hinton GE, Zemel RS, 1993. Autoencoders, minimum description length and helmholtz free energy. Proceedings of the 6th International Conference on Neural Information Processing Systems, San Francisco, CA, USA, p.3-10.
- Hinton GE, Dayan P, Frey BJ, et al., 1995. The "wake-sleep" algorithm for unsupervised neural networks. *Science*, 268(5214):1158-1161. <https://www.science.org/doi/abs/10.1126/science.7761831> <https://doi.org/10.1126/science.7761831>
- Ho J, Jain A, Abbeel P, 2020. Denoising diffusion probabilistic models, . <https://arxiv.org/abs/2006.11239>
- Hochreiter S, Schmidhuber J, 1997. Long short-term memory. *Neural computation*, .
- Huang G, Liu Z, Van Der Maaten L, et al., 2017. Densely connected convolutional networks. Proceedings of the IEEE conference on computer vision and pattern recognition, p.4700-4708.
- Hughes JF, van Dam A, McGuire M, et al., 2014. Computer Graphics: Principles and Practice, 3rd Edition. Addison-Wesley.
- , 2019. Automatic Machine Learning: Methods, Systems, Challenges. Springer.
- Hyvärinen A, Oja E, 1997. A fast fixed-point algorithm for independent component analysis. *Neural Computation*, 9:1483-1492.
- Hyvärinen A, 1997. A family of fixed-point algorithms for independent component analysis. *IEEE Int Conf on Acoustics, Speech and Signal Processing (ICASSP)*, p.3917-3920.
- Jin C, Netrapalli P, Jordan MI, 2019. What is local optimality in nonconvex-nonconcave minimax optimization? *preprint arXiv:190200618*, . <https://arxiv.org/abs/1902.00618>
- Jolliffe I, 1986. Principal Component Analysis. Springer-Verlag, New York, NY.
- Jonathan Bennett J, Carbery A, Christ M, et al., 2007. The brascamp-lieb inequalities: Finiteness, structure and extremals. *Geometric and Functional Analysis*, 17.
- Josselyn SA, Tonegawa S, 2020. Memory Engrams: Recalling the past and imagining the future. *Science*, 367.
- Kakade SM, 2001. A natural policy gradient. *Advances in Neural Information Processing Systems*, 14.
- Kanwisher N, McDermott J, Chun MM, 1997. The fusiform face area: A module in human extrastriate cortex specialized for face perception. *Journal of Neuroscience*, 17(11):4302-4311. <https://www.jneurosci.org/content/17/11/4302> <https://doi.org/10.1523/JNEUROSCI.17-11-04302.1997>

- Kanwisher N, 2010. Functional specificity in the human brain: A window into the functional architecture of the mind. *Proceedings of the National Academy of Sciences*, 107(25):11163-11170. <https://www.pnas.org/doi/abs/10.1073/pnas.1005062107> <https://doi.org/10.1073/pnas.1005062107>
- Keller GB, Mrsic-Flogel TD, 2018. Predictive processing: A canonical cortical computation. *Neuron*, 100(2):424-435.
- Kelley H, 1960. Gradient theory of optimal flight paths. *ARS Journal*, 30(10):947-954.
- Kingma DP, Welling M, 2013. Auto-encoding variational Bayes. *preprint arXiv:13126114*, .
- Kobyzev I, Prince SJ, Brubaker MA, 2021. Normalizing flows: An introduction and review of current methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(11):3964-3979. <https://doi.org/10.1109/tpami.2020.2992934>
- Koopman BO, 1931. Hamiltonian systems and transformation in Hilbert space. *Proceedings of the National Academy of Sciences of the United States of America*, 17 5:315-8.
- Kramer MA, 1991. Nonlinear principal component analysis using autoassociative neural networks. *Aiche Journal*, 37:233-243.
- Kriegeskorte N, Mur M, Ruff DA, et al., 2008. Matching categorical object representations in inferior temporal cortex of man and monkey. *Neuron*, 60(6):1126-1141.
- Krizhevsky A, Sutskever I, Hinton GE, 2012. Imagenet classification with deep convolutional neural networks. Advances in neural information processing systems, p.1097-1105.
- Kulkarni TD, Whitney WF, Kohli P, et al., 2015. Deep convolutional inverse graphics network. Advances in Neural Information Processing Systems, 28.
- LeCun Y, Browning J, 2022. What AI can tell us about intelligence. *NOËMA Magazine*, .
- LeCun Y, Bottou L, Bengio Y, et al., 1998. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, .
- LeCun Y, Bengio Y, Hinton G, 2015. Deep learning. *Nature*, 521(7553):436-444.
- LeCun Y, 2022. A path towards autonomous machine intelligence. *preprint posted on openreview*, .
- Lei N, Su K, Cui L, et al., 2017. A geometric view of optimal transportation and generative model. *preprint arXiv:171005488*, <https://arxiv.org/abs/1710.05488>
- Levoy M, Hanrahan P, 1996. Light field rendering. Proceedings of the 23rd annual conference on Computer graphics and interactive techniques, p.31-42.
- Li G, Wei Y, Chi Y, et al., 2020. Breaking the sample size barrier in model-based reinforcement learning with a generative model. Advances in Neural Information Processing Systems, 33:12861-12872.
- Ma Y, Soatto S, Kosecka J, et al., 2004. An Invitation to 3D Vision: From Images to Geometric Models. Springer-Verlag.
- Ma Y, Derksen H, Hong W, et al., 2007. Segmentation of multivariate mixed data via lossy data coding and compression. *IEEE Trans PAMI*, .
- MacDonald J, Wäldchen S, Hauch S, et al., 2019. A rate-distortion framework for explaining neural network decisions. *preprint arXiv:190511092*, .
- Marcus G, 2020. The next decade in AI: four steps towards robust artificial intelligence. *CoRR*, abs/2002.06177. <https://arxiv.org/abs/2002.06177>
- Marr D, 1982. Vision. MIT Press.
- Mayr O, 1970. The Origins of Feedback Control. MIT Press.
- McCloskey M, Cohen NJ, 1989. Catastrophic interference in connectionist networks: The sequential learning problem. 24:109-165.
- Mildenhall B, Srinivasan PP, Tancik M, et al., 2020. Nerf: Representing scenes as neural radiance fields for view synthesis. ECCV.
- Newell A, Simon HA, 1972. Human Problem Solving. Prentice Hall, New Jersey, USA.
- Ng A, Russell S, 2000. Algorithms for inverse reinforcement learning. *ICML '00 Proceedings of the Seventeenth International Conference on Machine Learning*, .
- Olshausen BA, Field DJ, 1996. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583):607.
- Oord Avd, Li Y, Vinyals O, 2018. Representation learning with contrastive predictive coding. *preprint arXiv:180703748*, .
- Osband I, Roy BV, 2014. Model-based reinforcement learning and the eluder dimension. Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 1, Cambridge, MA, USA, p.1466-1474.
- Pai D, Psenka M, Chiu CY, et al., 2022. Pursuit of a discriminative representation for multiple subspaces via sequential games. *preprint arXiv:220609120*, .
- Papayan V, Han X, Donoho DL, 2020. Prevalence of neural collapse during the terminal phase of deep learning training. *preprint arXiv:200808186*, .
- Patterson D, Gonzalez J, Hölzle U, et al., 2022. The carbon footprint of machine learning training will plateau, then shrink. *preprint arXiv:220405149*, <https://arxiv.org/abs/2204.05149>
- Quinlan JR, 1986. Induction of decision trees. *Machine Learning*, .
- Rao RPN, Ballard DH, 1999. Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. *Nat Neurosci*, 2(1):79-87.
- Rifai S, Vincent P, Muller X, et al., 2011. Contractive auto-encoders: Explicit invariance during feature extraction. In International Conference on Machine Learning, p.833-840.
- Rissanen J, 1989. Stochastic Complexity in Statistical Inquiry Theory. World Scientific Publishing Co., Inc., USA.
- Roberts DA, Yaida S, Hanin B, 2022. The Principles of Deep Learning Theory. Cambridge University Press (<https://deeplearningtheory.com>).
- Rosenblatt F, 1958. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65 6:386-408.
- Rumelhart DE, Hinton GE, Williams RJ, 1986. Learning representations by back-propagating errors. *Nature*, .
- Russell S, Norvig P, 2020. Artificial Intelligence: A Modern Approach. Pearson.
- Sastry S, 1999. Nonlinear Systems: Analysis, Stability, and Control. Springer.
- Shannon CE, 1948. A mathematical theory of communication. *The Bell System Technical Journal*, 27(3):379-423. <https://doi.org/10.1002/j.1538-7305.1948.tb01338.x>
- Shazeer N, Mirhoseini A, Maziarz K, et al., 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. ICLR, <https://openreview.net/pdf?id=BlckMDqIq>
- Shum HY, Chan S, Kang S, 2011. Image-Based Rendering. Springer US, <https://books.google.com/books?id=GeObcQAACAAJ>
- Silver D, Huang A, Maddison CJ, et al., 2016. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484-489.

- Silver D, Schrittwieser J, Simonyan K, et al., 2017. Mastering the game of go without human knowledge. *nature*, 550(7676):354-359.
- Simon HA, 1969. The Sciences of the Artificial. MIT Press, Cambridge, MA, USA.
- Srivastava A, Valkoz L, Russell C, et al., 2017. VeeGAN: Reducing mode collapse in GANs using implicit variational learning. *Advances in Neural Information Processing Systems*, p.3310-3320.
- Srivastava RK, Greff K, Schmidhuber J, 2015. Highway networks. *preprint arXiv:150500387*, .
- Sutton R, Barto A, 2018. Reinforcement Learning: An Introduction. MIT Press.
- Szegedy C, Zaremba W, Sutskever I, et al., 2013. Intriguing properties of neural networks. *arXiv:13126199*, .
- Szeliski R, 2022. Computer Vision: Algorithms and Applications, 2nd Edition. Springer-Verlag.
- Tenenbaum JB, de Silva V, Langford JC, 2000. A global geometric framework for nonlinear dimensionality reduction. *Science*, 290(5500):2319-2323.
- Tishby N, Zaslavsky N, 2015. Deep learning and the information bottleneck principle. *IEEE Information Theory Workshop*.
- Tong S, Dai X, Wu Z, et al., 2022. Incremental learning of structured memory via closed-loop transcription. *preprint arXiv:220205411*, .
- Uehara M, Zhang X, Sun W, 2021. Representation learning for online and offline RL in low-rank MDPs. *arXiv preprint arXiv:211004652*, .
- Vaswani A, Shazeer N, Parmar N, et al., 2017. Attention is all you need. *preprint arXiv:170603762*, .
<https://arxiv.org/abs/1706.03762>
- Vidal R, Ma Y, Sastry S, 2016. Generalized Principal Component Analysis. Springer Verlag.
- Vidal R, 2022. Attention: Self-expression is all you need, .
<https://openreview.net/forum?id=MmujBCLawFo>
- Vinyals O, Babuschkin I, Czarnecki WM, et al., 2019. Grandmaster level in StarCraft II using multi-agent reinforcement learning. *Nature*, 575(7782):350-354.
- Wang T, Buchanan S, Gilboa D, et al., 2021. Deep networks provably classify data on curves. *Advances in Neural Information Processing Systems*, 34.
- Wiatowski T, Bölcskei H, 2018. A mathematical theory of deep convolutional neural networks for feature extraction. *IEEE Transactions on Information Theory*, .
- Wiener N, 1948. Cybernetics. MIT Press, Cambridge, Mass.
- Wiener N, 1961. Cybernetics, 2nd edition. MIT Press, Cambridge, Mass.
- Wisdom S, Powers T, Pitton J, et al., 2017. Building recurrent networks by unfolding iterative thresholding for sequential sparse recovery. 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), p.4346-4350.
- Wood E, Baltrušaitis T, Hewitt C, et al., 2021. Fake it till you make it: Face analysis in the wild using synthetic data alone. *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, p.3681-3691.
- Wright J, Ma Y, 2022. High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications. Cambridge University Press (<https://book-wright-ma.github.io>).
- Wright J, Tao Y, Lin Z, et al., 2008. Classification via minimum incremental coding length (MICL). *Advances in Neural Information Processing Systems*, p.1633-1640.
- Xie S, Girshick R, Dollár P, et al., 2017. Aggregated residual transformations for deep neural networks. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), p.5987-5995.
- Yang Z, Yu Y, You C, et al., 2020. Rethinking bias-variance trade-off for generalization of neural networks. *International Conference on Machine Learning*.
- Yildirim I, Belledonne M, Freiwald W, et al., 2020. Efficient inverse graphics in biological face processing. *Science advances*, 6(10):eaax5979.
- Yu A, Fridovich-Keil S, Tancik M, et al., 2021. Plenoxels: Radiance fields without neural networks. *preprint arXiv:211205131*, .
- Yu Y, Chan KHR, You C, et al., 2020. Learning diverse and discriminative representations via the principle of maximal coding rate reduction. *Advances in Neural Information Processing Systems*, 33:9422-9434.
- Zhai Y, Yang Z, Liao Z, et al., 2020. Complete dictionary learning via ℓ^4 -norm maximization over the orthogonal group. *Journal of Machine Learning Research*, 21(165):1-68.
- Zhu JY, Park T, Isola P, et al., 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. *Proceedings of the IEEE international conference on computer vision*, p.2223-2232.
- Zoph B, Le QV, 2017. Neural architecture search with reinforcement learning,
<https://arxiv.org/abs/1611.01578>