

Toward an Over-parameterized Direct-Fit Model of Visual Perception

Xin Li

xin.li@ieee.org

Abstract

In this paper, we revisit the problem of computational modeling of simple and complex cells for an over-parameterized and direct-fit model of visual perception. Unlike conventional wisdom, we highlight the difference in parallel and sequential binding mechanisms between simple and complex cells. A new proposal for abstracting them into space partitioning and composition is developed as the foundation of our new hierarchical construction. Our construction can be interpreted as a product topology-based generalization of the existing k-d tree, making it suitable for brute-force direct-fit in a high-dimensional space. The constructed model has been applied to several classical experiments in neuroscience and psychology. We provide an anti-sparse coding interpretation of the constructed vision model and show how it leads to a dynamic programming (DP)-like approximate nearest-neighbor search based on ℓ_∞ -optimization. We also briefly discuss two possible implementations based on asymmetrical (decoder matters more) auto-encoder and spiking neural networks (SNN), respectively.

I. INTRODUCTION

How do we learn to see in the first six months after birth? To answer this question, David Hubel and Torsten Wiesel conducted pioneering experiments in the 1950s, leading to the discovery of simple and complex cells [43]. Inspired by their discovery, David Marr developed a theory of the neocortex [62] in 1970 and a theory of the hippocampus [63] in 1971. His computational investigation of vision [61] was published in 1982 after his death. The construction of neocognitron by Fukushima [31] and connectionist models by LeCun [50] in the 1980s represented the continuing effort to construct biologically plausible computational models for visual perception. Wavelet theory [57, 23] and sparse coding [70, 71] in the 1990s further supplied mathematical formulations of multi-resolution analysis for scale-invariant representation of images. Rapid advances in deep learning [33, 51], especially the class of over-parameterized models [7, 6] have expedited both the theory and practice of data-driven/learning-based visual processing.

Despite the great progress of today, the gap between biological and artificial vision remains significant in the following aspects. First, the network architecture of the convolutional neural network (CNN) is characterized by the pooling of layers, which reduces the dimensionality of the input data. This is in sharp contrast to the increase in the number of neurons and synapses as we move from the lower layer (e.g. V1) to the higher layer (e.g. V4) of neocortex. This anatomical finding has inspired H. Barlow to revise his redundancy reduction hypothesis into the redundancy exploitation hypothesis [9] in 2001. Second, although recurrent neural networks including long short-term memory (LSTM) [42] take into account temporal dynamics and have found successful applications in 1D signal analysis, the role of memory in visual perception has remained explored. In the human vision system, the hippocampus is known to play a critical role in various cognition tasks, including memory consolidation and novelty detection [49]. Finally, it remains a mystery how the human brain can manage to achieve the objectives of learning and memory with more than 100-1000 trillion synapses and a power budget of less than 20W. The challenge of breaking the conventional von Neumann architecture built upon the Turing machine remains a holy grail in neuromorphic computing.

The motivation behind this paper is two-fold. On the one hand, both the human brain and CNN are characterized by the ability to optimize an astronomical number of synaptic weights [20]. The class of over-parameterized models [7, 6] has shown some counterintuitive properties, such as double descent [68]. Analytical tools such as neural tangent kernel (NTK) offer an approach to understanding over-parameterization in the Hilbert space, but, like all kernel methods, they are not compatible with the

recursion strategy (e.g. dynamic programming that builds upon the optimality of substructures [13]). We seek to understand overparameterization in the framework of optimizing hierarchical representations [85]. On the other hand, an evolutionary perspective on biological and artificial neural networks [34] offers a direct-fit approach in brute force. Such a deceptively simple model, when combined with overparameterized optimization, offers an appealing solution to increase the generalization (predictive) power without explicitly modeling the unknown generative structure underlying sensory inputs.

In this paper, we construct an over-parameterized direct-fit model for visual perception. Unlike the conventional wisdom of abstracting simple and complex cells, we use space partitioning and composition as the building block of our hierarchical construction. In addition to biological plausibility, we offer a geometric analysis of our construction in topological space (i.e., topological manifolds without the definition of a distance metric or an inner product). Our construction can be interpreted as a product-topology-based generalization of the existing k-d tree [14], making it suitable for brute-force direct-fit in a high-dimensional space. In the presence of novelty/anomaly, a surrogate model that mimics the escape mechanism of the hippocampus can be activated for unsupervised continual learning [98]. The constructed model has been applied to several classical experiments in neuroscience and psychology. We also provide an anti-sparse coding interpretation [46] of the constructed vision model and present a dynamic programming (DP)-like solution to approximate nearest neighbor in high-dimensional space. Finally, we briefly discuss two possible network implementations of the proposed model based on asymmetric autoencoder [69] and spiking neural networks (SNN) [45], respectively.

II. NEUROSCIENCE FOUNDATION

A. Dichotomy: Excitatory and Inhibitory Neurons

During the work of Wilson and Cowan in the 1970s [92, 93], they made the crucial assumption that “all nervous processes of any complexity depend on the interaction of excitatory and inhibitory cells.” Using phase plane methods, they have shown simple and multiple hysteresis phenomena and limit cycle activity with localized populations of model neurons. Their results, more or less, offer the primitive basis for memory storage, namely stimulus intensity, which can be coded in both the average spike frequency and the frequency of periodic variations in the average spike frequency [78]. However, such ad hoc sensory encoding cannot explain the sophistication of learning, memory, and recognition associated with higher functions.

B. Hebbian Learning and Anti-Hebbian Learning

Hebbian learning [38] is a dogma that claims that an increase in synaptic efficacy arises from repeated and persistent stimulation of a presynaptic cell by a postsynaptic cell. Hebbian learning rule is often summarized as “cells that fire together wire together”. The physical implementation of the Hebbian learning rule has been well studied in the literature, for example, through spike timing-dependent plasticity (STDP) [16]. The mechanism of STDP is to adjust the connection strengths based on the relative timing of some neuron’s input and output action potentials. STDP as a Hebbian synaptic learning rule has been demonstrated in various neural circuits, from insects to humans.

By analogy to excitatory and inhibitory neurons, it has been suggested that a reversal of Hebb’s postulate, named anti-Hebbian learning, dictates the reduction (rather than increase) of the synaptic connectivity strength between neurons following a firing scenario. Synaptic plasticity that operates under the control of an anti-Hebbian learning rule has been found to occur in the cerebellum [12]. More importantly, local anti-Hebbian learning has been shown to be the foundation for forming sparse representations [27]. By connecting a layer of simple Hebbian units with modifiable anti-Hebbian feedback connections, one can learn to encode a set of patterns into a sparse representation in which statistical dependency between the elements is reduced while preserving the information. However, the sparse coding represents only a local approximation of the sensory processing machinery. To extend it to global (nonlocal) integration, we have to assume an additional postulate, called the “hierarchical organization principle”, which we will introduce in the next section.

C. Simple and Complex Cells in V1

These two classes of cells were discovered by Torsten Wiesel and David Hubel in the early 1960s [43]. Simple cells respond primarily to oriented edges and gratings, which can be mathematically characterized by Gabor filters [24]. Complex cells also respond to oriented structures; unlike simple cells, they have a degree of spatial invariance. The difference between receptive fields and the characteristics of simple and complex cells has inspired the invention of the neocognitron by Fukushima in 1979 [31], which foresaw the subsequent convolutional neural network. The hierarchical convergent nature of visual processing has also inspired the construction of the HMAX model in 1999 [81].

An important observation with the difference between simple and complex cells, as described in [43], is their neural circuits and the corresponding temporal dynamics. Simple cells are built from center-surrounding cells, which require *simultaneous summation*. On the contrary, activation of the complex cell by a moving stimulus requires *successive activation* of many simple cells. Therefore, the spatial invariance of complex cells is achieved by summation and integration of the receptive fields of simple cells. Mathematical modeling of complex cells has been extensively studied in the literature (e.g., energy model [3]). However, the abstraction strategies taken in this paper will be different from those in the open literature.

D. Mountcastle's Universal Principle

In 1978 Mountcastle suggested a universal processing principle that has been acclaimed as the Rosetta Stone of neuroscience. According to Mountcastle, all parts of the neocortex operate according to a common principle, the cortical column being the unit of computation [67]. If Mountcastle were correct, the “simple discovery” made by Hubel and Wiesel might have deeper implications in the mechanism of visual processing beyond V1. Along this line of reasoning, the striking difference between simultaneous and successive activation of simple and complex cells might illustrate a fundamental contrast between two classes of binding mechanism among neurons.

In visual perception, it has been hypothesized that the characteristics of individual objects are bound / segregated by Hebbian / anti-Hebbian learning of different groups of neurons [66]. We conjecture that there exist two types of binding mechanism (parallel vs. sequential) that are analogous to combinatorial and sequential logic in digital circuits. The former plays the role of integrating spatially overlapped parts into a whole (e.g., a horizontal edge and a vertical edge form a letter “T”) or multiple features of the same object into a coherent perception (e.g., the age and gender of a face) in object recognition. Parallel binding can be interpreted as an extension of von der Malsburg's correlation theory [58]. This is the mechanism adopted by the dorsal stream to support the task of object vision. The latter is at the core of integrating spatially non-overlapped parts into a whole (e.g., the concatenation of letters into a word) in spatial vision, which belongs to the ventral stream/pathway. Sequential binding is closely related to the formation of short-term memory (e.g., Miller's law [65]) and long-term memory (e.g., Atkinson–Shiffrin model [8]) in the brain. The fundamental difference between parallel and sequential binding is that the former is invariant to permutation (the ordering of parts does not affect the perception of the whole), while the latter is sensitive to the ordering of the neuronal groups.

Binding by neuronal synchrony has been widely recognized in the literature; however, the binding problem is often thought to suffer from the so-called “superposition catastrophe” [91]. The combination coding argument often faces the dilemma of the curse of dimensionality, and it has been suggested that a representation with hierarchical structure can at least partially overcome this barrier of impracticality with combination coding. More importantly, we argue that our intuition about the capacity of the cortex might be misleading and our understanding of the power of hierarchical structures is inadequate [35]. If the curse of dimensionality can become a blessing [18], the combination coding can be made compatible with the hypothesis of redundancy exploitation [9].

III. CONSTRUCTION OF AN OVER-PARAMETERIZED DIRECT-FIT MODEL

In neuroscience, the principle of hierarchical organization can be roughly stated as follows: The nested structure of the physical world is mirrored by the hierarchical organization of the neocortex [35]. Unlike the mathematical construction of wavelets [23, 57], we envision that nature has discovered an elegant “elementary” solution in topological space (without the extra structure of distance metric or inner product) to manage the complexity of sensory stimuli in the physical world. We propose to study the following problem as the fruit-fly problem in visual perception [39].

Problem Formulation of Visual Perception

Given a visual stimulus (e.g., a sequence of images) as input, group/cluster them into different classes in an unsupervised manner.

The solution, as manifested by an infant’s development of the visual cortex (primarily for ventral stream for object vision) during the first six months after birth, lies in a novel construction of a hierarchical direct-fit model based on simple and complex cells. As argued by Jean Piaget [76], the ordering of mathematical spaces in early children cognitive development (topology before geometry) is the opposite to that of what we have learned in school (topology after geometry). Therefore, we attempt to construct our visual perception model in topological space with the least amount of assumed mathematical structures.

A. Preliminary on Topological Space

To facilitate our abstraction of simple and complex cells by subspace and product topology, we briefly review the basics of topological space as follows. We will follow the axiomatization of Felix Hausdorff to construct the topological space using neighborhood as the building block. Let \mathcal{N} denote the neighborhood function assigning to each point $x \in \mathbf{X}$ a non-empty subset $\mathcal{N}(x) \subseteq \mathbf{X}$. Then the following axioms must be satisfied for \mathbf{X} with \mathcal{N} to be called a topological space.

- 1) If \mathbf{N} is a neighborhood of x (i.e., $\mathbf{N} \in \mathcal{N}(x)$), then $x \in \mathbf{N}$;
- 2) If \mathbf{N} is a subset of \mathbf{X} and includes a neighborhood of x , then \mathbf{N} is a neighborhood of x ;
- 3) The intersection of two neighborhoods of x is a neighborhood of x ;
- 4) Any neighborhood \mathbf{N} of x includes a neighborhood \mathbf{M} of x such that \mathbf{N} is a neighborhood of each point in \mathbf{M} .

Note that the fourth axiom plays the role of linking the neighborhoods of different points in \mathbf{X} together. Since no distance metric is defined, we need to define the basis as the starting point for defining a topology.

Basis of a Topology: Let \mathbf{X} be a set, and suppose that \mathcal{B} is a collection of subsets of \mathbf{X} . Then \mathcal{B} is a basis for some topology in \mathbf{X} if and only if the following two conditions are satisfied: a) $\cup_{\mathbf{B} \in \mathcal{B}} \mathbf{B} = \mathbf{X}$; b) If $\mathbf{B}_1, \mathbf{B}_2 \in \mathcal{B}$ and $x \in \mathbf{B}_1 \cap \mathbf{B}_2$, there exists an element $\mathbf{B}_3 \in \mathcal{B}$ such that $x \in \mathbf{B}_3 \subseteq \mathbf{B}_1 \cap \mathbf{B}_2$.

In this work, we will only consider the neighborhood basis, which is defined by

Neighborhood Basis: A neighborhood basis is a subset $\mathcal{B} \subseteq \mathcal{N}(x)$ such that for all $\mathbf{V} \in \mathcal{N}(x)$, there exists some $\mathbf{B} \in \mathcal{B}$ such that $\mathbf{B} \subseteq \mathbf{V}$. In other words, for any neighborhood \mathbf{V} we can find a neighborhood \mathbf{B} on the basis of the neighborhood contained in \mathbf{V} .

With the above setup, the objective is to construct a hierarchical direct-fit model in the Hausdorff space (a.k.a. topological manifold), which generalizes the existing multi-resolution analysis in the Hilbert space [57]. Following our intuition above, simple and complex cells will be abstracted into subspace and product topology [54], respectively. Formally, we have the following.

Subspace Topology: Let \mathbf{X} be a topological space and let $\mathbf{S} \subseteq \mathbf{X}$ be any subset. Then $\mathcal{T}_S = \{\mathbf{U} \subseteq \mathbf{S} : \mathbf{U} = \mathbf{S} \cap \mathbf{V} \text{ for some open subset } \mathbf{V} \subseteq \mathbf{X}\}$ is the subspace topology.

Product Topology: Suppose that $\mathbf{X}_1, \dots, \mathbf{X}_n$ are arbitrary topological spaces. In its Cartesian product $\mathbf{X}_1 \times \dots \times \mathbf{X}_n$, the product topology is generated on the following basis: $\mathcal{B} = \{\mathbf{U}_1 \times \dots \times \mathbf{U}_n : \mathbf{U}_i \text{ is an open subset of } \mathbf{X}_i, i = 1, \dots, n\}$.

Both subspace and product topologies have their uniqueness in terms of satisfying the characteristic property [48]. The geometric intuition behind our construction of the new hierarchical model is best illustrated by the duality between space partitioning (i.e., subspace topology) and composition (i.e., product topology).

B. Computational Modeling of Neocortex

Simple Cells play the role of space partitioning, which can be abstracted as subspace topology [48]. A good proxy model to study the concept of space partitioning is the k-dimensional tree (k-d tree) [14], a space partitioning data structure to organize points in a k-dimensional space. The k-d tree structure can be interpreted as a class of binary space partitioning trees that extends the binary search tree (BST) for sorted arrays. It directly fits the data using hyperplanes to recursively partition the k-dimensional space. A simple variant of the k-d tree, named the random projection tree (rp tree) [22], is capable of automatically adapting to the low-dimensional structure of the data without explicit manifold learning.

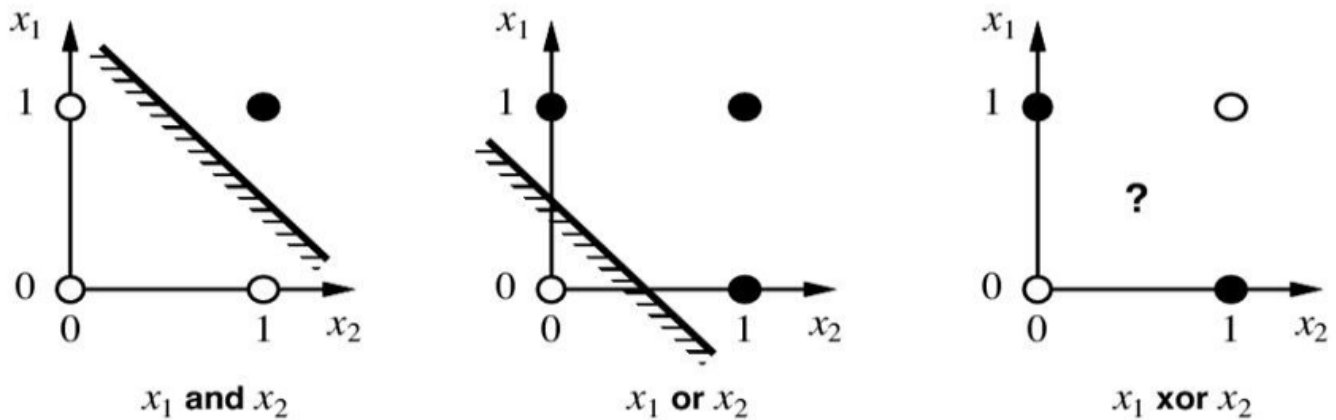


Fig. 1: XOR is not linearly separable but can be decomposed into the linear difference between linearly separable sets (the base case for proving the convex decomposition lemma).

The limitations of linear separability (e.g., single-layer perceptron [82]) are well known. Conventional wisdom of addressing these limitations is to introduce a nonlinear activation unit or a hidden layer (e.g., multi-layer perceptron). However, we note that some linearly non-separable set (e.g., the well-known XOR as shown in Fig. 1) can be decomposed into the linear difference between two linearly separable sets ($X \oplus Y = (X|Y) - (X \wedge Y)$). Such an observation, when combined with the data structure such as k-d/rp trees, offers a refreshing perspective of adapting to the intrinsic low-dimensional structure in high-dimensional data. That is, instead of nonlinear dimensionality reduction (e.g., the well-known Johnson-Lindenstrauss lemma [89]), we can use only a linear combination of space partitioning and differencing to approximate an arbitrary nonconvex object based on the following lemma.

Lemma 1: Convex Decomposition Lemma for Simple Cells

Any nonconvex object X can be decomposed into a finite set of differences between convex objects.

Sketch of the Proof Prove by induction with n , which is the minimum number of convex objects $\{X_1, \dots, X_n\}$ for the generation of X . Starting with $n = 2$, it is easy to construct the solution for $X = X_1 \pm X_2$. Suppose that the statement is true for the case of $n = k$, then the solution to $n = k + 1$ can be constructed similarly to $n = 2$.

Note that k-d trees have the nice property of facilitating NN/kNN search in metric space [14, 29]; but its performance degrades in high-dimensional space due to the curse of dimensionality. More importantly, one of the surprising behavior of distance metrics in high-dimensional spaces [4] is that ℓ_p ($p \rightarrow 0$) behaves much better than the popular Euclidean distance metric. Counterintuitively, in a high-dimensional space, proximity-based queries such as the NN search are meaningless and unstable because the discrimination between the nearest and farthest neighbor becomes poor. As rigorously shown in [4], the relative contrast provided by a norm with a smaller parameter p is more likely to dominate another norm with a larger parameter p as the dimensionality increases. The fractional distance concentration [28] dictates that we either work with topological space (instead of metric space) or use ℓ_0 as a pseudo-norm that recognizes

the limitations of distance metrics. Interestingly, the problem of ℓ_∞ -optimization has shown to be useful for an approximate NN search with anti-sparse coding [46].

Complex Cells play the role of the space composition, which can be abstracted by product topology [48]. Note that our intuition is consistent with the max/sum-pooling operation in HMAX model [81] because the objective is to achieve spatial invariance within an increased field of view. The difference lies in the way of abstraction - from simple to complex cells, we ask what will be its dual operation of space partitioning. Along this line of reasoning, if simple cells are responsible for linear separability without the change of dimensionality; complex cell must be able to increase the dimensionality for handling more sophisticated objects.

A common criticism of increasing dimensionality is concern about the so-called *curse of dimensionality* [13]. There are several ways to address this problem. First, recently proposed overparameterized models [7] or direct-fit models [34] suggest that dimensionality can be a blessing when a large amount of training data is available due to a counterintuitive phenomenon called “concentration of measure” [52]. Second, experimental studies such as [18] have demonstrated the blessing of dimensionality in face verification applications, which is consistent with biological findings [64]. More importantly, the low-dimensional manifold structure can still be preserved after space composition because of the following lemma.

Lemma 2: Product-Manifold Lemma [54] for Complex Cells

In topology, a topological manifold is a topological space that locally resembles the real n -dimensional Euclidean space. Let M be a topological m -manifold and N be a topological n -manifold, then $M \times N$ (Cartesian product of M and N) is a topological $(m + n)$ -manifold.

This lemma explains the blessing of dimensionality in that the manifold structure is easier to discover in a high-dimensional space. Note that manifold learning in a higher-dimensional space requires more training data. For example, the Cartesian product of a horizontal edge and a vertical edge will produce several combinations including “T”, “+”, “┌”, “┐”, and “└”. Through combination coding, our direct-fit model stores different patterns like k-d tree (i.e., the centroids of vector quantization or dictionary atoms of sparse coding) but the combination of those patterns will be further enumerated through product topology to support the direct-fit at the higher hierarchy. This enumeration perspective differs from our approach from the HMAX model [81] in which no discrimination was considered for the combination of basic patterns. We argue that the combinatorial coding argument is consistent with recently developed direct-fit model [34], as we will elaborate next.

C. Hierarchical Direct-Fit Model with Manifold-based Novelty Detection

Hierarchical Model Construction combines layers of simple and complex cells such as HMAX [81] or neocognitron [31], but with an important distinction. The network architecture of our model is not convergent, but *divergent* - one way of generalizing is to still use pooling layers, but we will consider many parallel pooling layers simultaneously. This divergent architecture is inspired by the hypothesis of redundancy exploitation [9] advocated by H. Barlow, since the number of neurons does not decrease, but increases significantly as we move to the higher level of the neocortex. It is natural to ask why higher functions in visual recognition need more neurons. In addition to the argument with combination coding, we note that achieving spatial invariance by max-pooling loses the resolution. To compensate for this sacrifice, context aggregation by dilated convolutions [97] has been developed for semantic segmentation. Mathematically, a dilated convolution is equivalent to a convolution without the follow-up max-pooling operation. Alternatively, we can still use max-pooling, but consider the generalization of dilated convolutions to deformable convolutions [21]. To achieve invariance to local geometric transformations, we can consider a pool of deformable models instead of a single one. Such a combination of space composition with deformable kernels allows us to generate invariant representations with product topology in a bottom-up fashion.

Based on the constructed hierarchical model, we note that the combination of space partitioning (by simple cells) and space composition (by complex cells) can be recursively applied to obtain an

over-parameterized model in higher-dimensional space. This recursion is conceptually similar to multi-resolution analysis [57] but operates in a data-adaptive manner (note that we have given up the structures of basis function and inner product in the Hilbert space). At each level, the concatenation of simple and complex cells will map the visual stimuli onto a sequence of invariant representations with increasing dimensionality (field of view). Due to the product-manifold lemma, the number of training data required by our direct-fit model will not exponentially increase (thus, avoiding the curse of dimensionality). Instead, a hierarchical model allows us to automatically adapt to low-dimensional structures in high-dimensional space by extending the k-d tree as follows.

Product Manifold Tree/Forest.

The product manifold tree (pm tree) can be defined as the dual operation of the classic k-d/rp tree. Instead of space partitioning, we recursively merge low-dimensional manifolds in subspaces into higher-dimensional manifolds through product topology. A pm forest uses a pm tree as its building block and consists of an ensemble of pm trees.

How to directly fit the data to the pm tree? Such a problem has been formulated as unsupervised learning in the literature of ML [41, 10]. Unfortunately, all existing work on unsupervised learning and clustering analysis assumes a fixed dimensionality with the input data and focuses on the learning of a distance metric. Our construction of the pm tree attempts to overcome such a barrier by generalizing k-d tree-based clustering with product topology. The performance of nearest-neighbor (NN)-search is known to degrade in high-dimensional space partially due to the surprising behavior of distance metric as dimensionality increases [4]. However, such limitation can be overcome by using the following lemma (note that we do not assume any definition of a distance metric in a topological manifold).

Lemma 3: Unsupervised Clustering Lemma on PM Tree.

Let $\mathbf{Z} = \mathbf{X} \times \mathbf{Y}$ denote the Cartesian product of two topological manifolds. For a vector $z = [xy]$, its neighborhood search can be carried out by taking the intersection of neighborhoods of $x \in \mathbf{X}$ and $y \in \mathbf{Y}$, respectively.

Sketch of the Proof. It is known that the subspaces and products of the Hausdorff spaces are Hausdorff. One property of product topology is that if \mathcal{B}_i is a basis for the topology of \mathbf{X}_i , then the set $\{\mathbf{B}_1 \times \dots \times \mathbf{B}_n : \mathbf{B}_i \in \mathcal{B}_i\}$ is a basis for the product topology on $\mathbf{X}_1 \times \dots \times \mathbf{X}_n$ [54]. The result follows directly from the definition of a neighborhood basis in the topological manifold.

Such hierarchical construction with interlaced simple and complex cell layers allows an organism to directly fit the visual stimuli into the hierarchical model in a bottom-up fashion. However, such a feedforward process alone is insufficient; it requires a control mechanism (negative feedback) for stability and an escape mechanism (novelty detection) for adaptation. The hippocampus plays an important role in memory consolidation (both old and new) and novelty detection [49]. One way of abstracting the functionality of the hippocampus in visual perception is that it serves the dual role: 1) a control mechanism for implementing negative feedback control loops; 2) an escape mechanism for out-of-distribution (OOD) samples. The first role has been extensively studied in the literature on deep learning (for example, the backpropagation algorithm [83]). Computational modeling of novelty detection includes statistical approaches [59] and neural network-based approaches [60]. Under the framework of our direct-fit hierarchical model, we can revisit these two mechanisms as follows.

Manifold-based Novelty Detection. In previous work [77], novelty detection was formulated as a twist on the manifold learning problem. First, it linearizes the parameterized manifold, capturing the underlying structure of the inlier distribution. The novelty score is then obtained by factoring in the probability and calculated with respect to the local coordinates of the tangent space of the manifold. However, the task of manifold learning is undertaken by an adversarial auto-encoder, which makes strong assumptions about the input images (e.g., fixed dimensions and inlier/outlier categories). By replacing the adversarial autoencoder with our newly constructed hierarchical model, we can achieve an improved generalization property in the following aspects. First, our construction of a product-manifold tree belongs to unsupervised learning because it directly fits the input data to a pre-constructed data structure. Second, similar to the k-d tree but working toward the opposite direction (increasing dimensionality), the product manifold tree can

dynamically accommodate novel events through continual learning. The escape mechanism allows us to incorporate a novel object into the tree so that it is treated as normal in the future.

IV. SPARSE CODING INTERPRETATION: DP-LIKE SOLUTION TO ℓ_∞ -OPTIMIZATION AND APPROXIMATE NEAREST NEIGHBOR SEARCH

A. Hierarchical Convolutional Sparse Coding

It is well known that sparse coding offers a powerful analysis of the mechanism of V1 [70, 71]. Meanwhile, sparse representations have also found promising applications in unsupervised learning, such as K-SVD dictionary learning [5] and non-negative matrix factorization [53]. Unlike predetermined dictionaries (e.g., wavelet [57]), data-adaptive dictionary learning such as K-SVD based sparseland [5] has led to a multi-layer formulation of convolutional sparse coding (ML-CSC) [73, 74], which provided an attractive new theoretical framework for analyzing CNN.

Multi-Layer Convolutional Sparse Coding (ML-CSC).

The new insight brought about by ML-CSC [73, 74] is to generalize the original sparse coding in a hierarchical (multi-layer) manner. Specifically, a multilayer convolutional sparse coding (ML-CSC) model can be constructed as follows. Suppose that \mathbf{X} is the input signal and a set of dictionaries is given by $\{\mathbf{D}_k\}_{k=1}^K$ where \mathbf{D}_k denotes the dictionary at the level k . Then an ML-CSC model can be written as: $\mathbf{X} = \mathbf{D}_1\mathbf{\Gamma}_1$, $\mathbf{\Gamma}_1 = \mathbf{D}_2\mathbf{\Gamma}_2$, ..., $\mathbf{\Gamma}_{K-1} = \mathbf{D}_K\mathbf{\Gamma}_K$ where $\mathbf{\Gamma}_i = [\mathbf{w}_1, \dots, \mathbf{w}_k]$ denotes the sparse coefficients at the level i . Following the convex approximation of ℓ_0 -optimization in [73], a layered thresholding algorithm runs recursively as follows.

$$\mathbf{\Gamma}_k = \mathcal{P}_{\beta_k}(\mathbf{D}_k^T \mathbf{\Gamma}_{k-1}), (k = 1, 2, \dots, K) \quad (1)$$

where \mathcal{P}_{β_k} is the standard thresholding operator and $\{\beta_k\}_{k=1}^K$ is the set of thresholds at level k .

As shown in [74], the ML-CSC model manages to decompose a signal $\mathbf{X} \in R^N$ into the superposition of multiple dictionaries $\mathbf{D} = \mathbf{D}_1\mathbf{D}_2\dots\mathbf{D}_K$ but the concatenation of these atoms, although it is overcomplete, remains in the space of the same dimension R^N . Unfortunately, both K-SVD and ML-CSC are still constructed within the Hilbert space without a change of dimensionality. To the best of our knowledge, no previous work exists in the open literature that extends data-adaptive sparse representation into varying dimensions. Note that wavelet theory achieves the objective of multi-resolution analysis by varying the support of basis functions. Our intuition is that a similar objective can be achieved for sparse data-adaptive representations if we can automatically adapt dictionary learning to the intrinsic low-dimensional structure of the data, such as rp trees [22]. A different way of generalizing the CSC is to construct a dictionary hierarchy with varying dimensionality.

Parallel to the product-manifold lemma, we are interested in developing a recursive strategy to decompose a high-dimensional sparse coding problem into the “product” of multiple lower-dimensional ones. Note that unlike existing work on wavelet decomposition dealing with basis construction in a Hilbert space, we opt to work with directly fitting the data, which is closer to the sparse PCA [99] and K-SVD [5]. Instead of formulating the dictionary learning in the Hilbert space with a fixed dimension (i.e., only vary dictionary size), we imagine a dynamic programming (DP)-like approach to dictionary learning from optimal substructures. The theoretical foundation of such a DP-like solution is the product manifold lemma; and sparse coding offers a well-established framework for exploiting the manifold constraint. We start from the following construction.

Hierarchical Convolutional Sparse Coding.

Let $\mathbf{X} = \mathbf{D}_x\mathbf{\Gamma}_x$ and $\mathbf{Y} = \mathbf{D}_y\mathbf{\Gamma}_y$ denote two dictionary coding schemes with $\mathbf{D}_x, \mathbf{D}_y \in R^{n \times m}$, ($m > n$). Then we start with a coding scheme in the direct-sum space R^{2n} by $\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} = \begin{bmatrix} \mathbf{D}_x & \mathbf{0} \\ \mathbf{0} & \mathbf{D}_y \end{bmatrix} \begin{bmatrix} \mathbf{\Gamma}_x \\ \mathbf{\Gamma}_y \end{bmatrix}$ and improve the sparsity by basis pursuit algorithm [19] using the composite dictionary $[\mathbf{D}_{x,i} \ \mathbf{D}_{y,j}]$, ($1 \leq i, j \leq m$) (total m^2 atoms). Such process of basis composition can be recursively applied to obtain sparse bases in higher-dimensional spaces (i.e., $R^n \rightarrow R^{2n} \rightarrow R^{4n} \dots$).

Our hierarchical extension of the CSC shares a spirit similar to that of dynamic programming in that optimal substructures (atoms) contribute to the (nearly) optimal solution (molecule). Mathematically, ℓ_0 -optimization is an NP-hard problem; but in biology, evolution does not have foresight and therefore the hierarchical systems generated by evolution might not have a globally optimal structure. What matters more appears to be the nearly decomposability of complex systems [85], which is closely related to the principle of dynamic programming (DP) [13]. Our intuition is that evolution does not need to pursue a globally optimal solution such as NN, but be satisfied with an approximate yet flexible solution so that the organism can adapt to the constantly evolving environment. Based on this observation, we connect the hierarchical CSC with a DP-like recursive solution to approximate NN (ANN) search next.

B. Anti-sparse coding for Approximate Nearest Neighbor (ANN) Search

Instead of ℓ_0 -optimization, we conjecture that ℓ_∞ -optimization (a.k.a., minimax optimization [36]) is a more appropriate framework for analyzing ventral stream processing for the following reasons. First, the strategy employed by V1 [70] follows Barlow's redundancy reduction principle [11]; however, Barlow has revised this principle to exploit redundancy in [9]. How to exploit the blessing of dimensionality by our hierarchical CSC calls for a fresh look at the definition of sparsity. Second, redundancy has been extensively exploited in information theory for reliable communication [84]. The neocortex faces a similar challenge of robustness to errors (e.g., sensory deprivation and lesions), especially for the high-level layers responsible for important decisions related to behavior. In the literature, it has been shown in [56, 30] that ℓ_∞ -optimization leads to the so-called Kashin representation (a.k.a. spread representation [30]) where all coefficients are of the same order of magnitude. Such a class of representations is known to robustly withstand errors in their coefficients. Third, the anti-sparse coding scheme based on ℓ_∞ -optimization is known to facilitate the search for the approximate nearest neighbor (ANN) [46]. Such an interesting connection implies that it is easier to construct a DP-like solution by decomposing the high-dimensional ANN search problem into several subproblems in projected subspaces.

Spread Representation For a given signal \mathbf{x} in n -dimensional space with $\mathbf{D}_x \in R^{n \times m}$, ($m > n$) representing the dictionary, we consider the following ℓ_∞ -optimization problem.

$$\alpha^* = \operatorname{argmin}_{\alpha} \|\alpha\|_{\infty}, \text{ s.t. } \mathbf{D}_x \alpha = \mathbf{x} \quad (2)$$

where $\|\alpha\|_{\infty} = \max |\alpha_i|, i \in \{1, \dots, m\}$. The solution to the above ℓ_∞ -optimization problem [30] boils down to a binarization scheme with components $m - n + 1$ reaching limits $\|\alpha\|_{\infty}$ and the remaining $n - 1$ between these two extreme values. This spread representation has led to an anti-sparse coding scheme for the ANN search [46].

Anti-sparse coding for ANN search. Given an inquiry vector \mathbf{y} , we can find its binarization $e(\mathbf{y}) = \operatorname{sign}(\alpha / \|\alpha\|_{\infty})$ by solving the anti-sparse coding problem.

$$J_{\lambda}(\alpha) = \frac{1}{2} \|\mathbf{D}_x \alpha - \mathbf{y}\|_2^2 + \lambda \|\alpha\|_{\infty}, \quad (3)$$

where λ is the Lagrangian multiplier. Then the problem of finding an ANN inquiry \mathbf{z} boils down to

$$\mathbf{NN}(e(\mathbf{z})) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} e(\mathbf{z})^T e(\mathbf{y}) \quad (4)$$

Combining anti-sparse coding with our hierarchical CSC, we can envision a DP-like recursive solution to the ANN search in high-dimensional space. Thanks to product sparse coding lemma, we can recursively construct a dictionary in a high-dimensional space from the direct product of dictionaries in low-dimensional spaces. It follows that an anti-sparse coding scheme in high-dimensional space can be obtained by decomposing it into subproblems in low-dimensional spaces (conceptually similar to the principle of pf dynamic programming). It should be noted that space partitioning, such as the k-d tree and

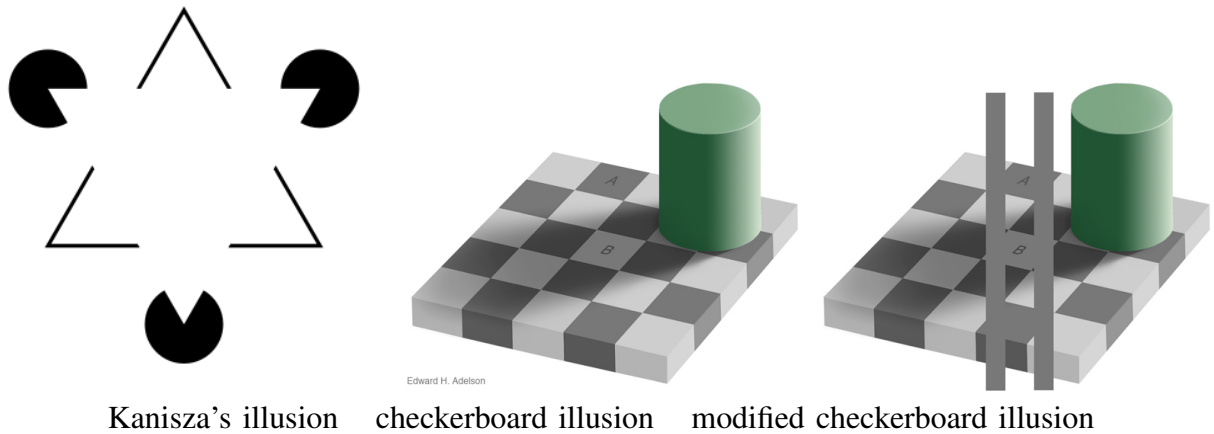


Fig. 2: Optical illusion of an illusory contour. The white triangle is perceived as an illusory contour (the illusion will disappear as the size of the triangle increases).

random projection [95] also supports the NN/ANN search [29]. Therefore, the strategy of ANN search might serve as a common currency to unify top-down and bottom-up processing mechanisms in visual perception.

C. Theoretical Analysis

Relative Invariance to Transformation and Illumination Variation. Unlike wavelet theory [57] in which scale/shift invariance is a property of basis functions, our direct-fit model achieves *relative* invariance to geometric transformations and changes in illumination by learning from the samples. Note that translational invariance is relatively easy to achieve with a maximum-pooling operator with an increased field of view. The invariance in the composition of geometric transformations is beyond the reach of conventional CNN. Our hierarchical model achieves such invariance through densely sampling the data space (guided by ego-motion). In other words, invariant representation arises from the interaction between sensory and motor system - as JJ Gibson comments on visual perception [32], “we see because we move; we move because we see.” Note that such invariance is relative - for rotation that is beyond a critical threshold, the performance of visual recognition degrades rapidly (escape mechanism will be activated).

Robustness to Occlusion and Corruption. The robustness of sparse coding recognition [94] is mainly attributed to the use of the ℓ_1 norm in convex optimization. The hierarchical extension further improves the robustness by exploiting the redundancy of the representation. As shown in [25], when an overcomplete system is *incoherent*, the optimally sparse approximation to noisy and noiseless data differs from at most a constant multiple of the noise level. Our intuition is that incoherence becomes easier to satisfy in a higher-dimensional space due to the phenomenon of concentration of measures [52]. In addition to the stability results proven in [73], additional robustness to missing and noisy observations comes from the ℓ_∞ -optimization as rigorously demonstrated by the uncertainty principle in [56].

Multiple Stable States. Our new minimax optimization perspective offers a refreshing interpretation about the bistability of anti-sparse coding. Note that \mathbf{x} and $-\mathbf{x}$ are always legitimate solutions to the same minimax optimization. This seems a salient difference from ℓ_p -norm in that ℓ_∞ -norm is the only one satisfying the symmetry constraint.

V. APPLICATIONS TO VISUAL PERCEPTION

A. Grandmother Cell Hypothesis and Face Identity Encoding

There is a long history of debate between localist and distributed representations in psychology and neuroscience [72]. In localist coding, a neuron codes for one familiar thing and does not directly contribute

to the representation of anything else [15]. In distributed representations, knowledge is coded in a distributed manner in the mind and brain. The discovery of grandmother cells (a.k.a. Jennifer Aniston cells) [79] can be interpreted as an extreme case of sparseness. The hypothesis of a grandmother cell is also closely related to the binding problem [91] that deals with the integration of individual features into a holistic experience. Our model is in support of the hypothesis regarding the grandmother cell. Following our analysis of the ML-CSC model, the sparsity is expected to increase as the depth of the network increases. This matches our intuition that more abstraction is obtained at higher levels, implying fewer nonzero coefficients required for a global sparse representation. The extreme case will be a single coefficient or neuron. However, it should be noted that there might be massive redundancy in coding any item, such as the grandmother's face. The locations of these redundant single-neuron encodings can be distributed in the cortex.

Meanwhile, we note that familiarity is an important constraint for the grandmother cell hypothesis. Familiarity and identity of faces are two related but different concepts encoded by the medial temporal lobe (MTL) and the fusiform face area [47]. Face familiarity is related to the bias of visual stimuli ; within the framework of sparse coding, faces within the social group (e.g., grandmother's face) represent the redundancy of training data. This redundancy corresponds to a denser sampling of a local region in the face space, which implies improved sparsity. Therefore, it is easier to observe the firing of single neurons in the class of familiar faces. On the contrary, unfamiliar faces (for example, the cross-race effect [96]) are often more difficult to recognize. This bias is associated with a degraded sparsity due to a poor sampling density. Therefore, it is more difficult to record the response of individual neurons, which is consistent with a recent finding on facial identity coding [17].

B. Optical Illusion and Context Adaptation

Optical illusions have been widely used by Gestalt psychologists as an experimental tool for studying perceptual organization in visual perception. For example, consider the famous Kanisza's illusion - an example of an illusory contour [90]. What has been less studied is the perturbed version of this illusion - i.e., as one increases the size of the black triangle and the distance among three pacmacs, it will become more and more difficult to perceive the illusory "white triangle" at the center. Such experimental findings cannot be explained by Gestalt theory because there is no prediction of the critical boundary condition for perceptual organization to fall apart. Our hierarchical direct-fit model can predict that such a threshold is determined by the size of the fovea centralis. When the distance is above this threshold, there is no previous experience (training data) to support the perceptual grouping of white triangles.

Another celebrated illusion related to context adaptation is Ted Adelson's checkerboard illusion [2] (see Fig. 2b). Two blocks with identical color values can appear to have opposite perceived luminance (brightness). Such an inconsistency is the consequence of perceptual organization; namely, the global constraint of perceiving a checkerboard dictates the interpretation at lower levels. By contrast, if two additional vertical strips having identical color are added, they interfere with the result of perceptual organization at the global level. Consequently, two blocks of the same color as two strips tend to be grouped together.

C. Multi-stable Perception and Novelty Detection

Multistable perception refers to the spontaneous alternation between two (or more) perceptual states when visual stimuli are inherently ambiguous (refer to Fig. 3). Multistable perception is a widely studied topic in visual perception [55] due to its close relationship to sensory awareness or consciousness. The neural basis of multi-stable perception [86] has emphasized the role of high-level brain mechanisms that are involved in actively selecting and interpreting sensory information from lower-level processes. However, no theoretical explanation about the deeper mechanism of multi-stable perception exists, not to mention the prediction about when this phenomenon will occur. Our hierarchical direct-fit model shed some insight to this long-standing open question in that bi-stable perception directly corresponds to x and

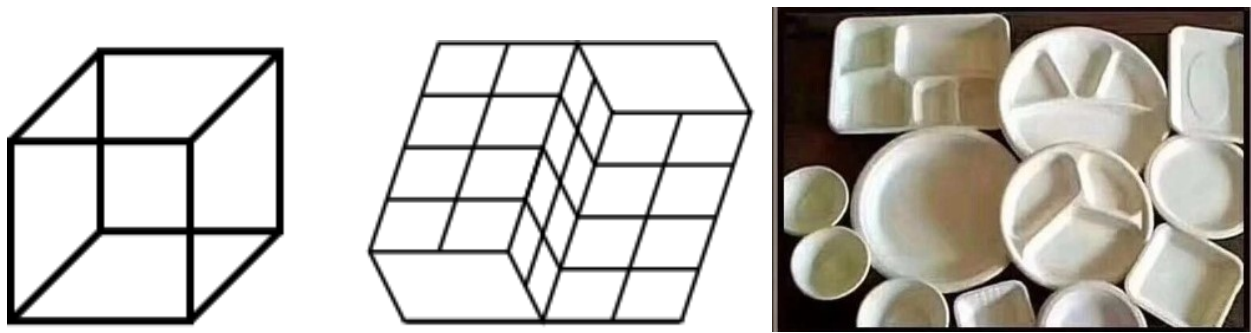


Fig. 3: Multi-stable perception (each image admits two competing interpretation results).

$-x$ in our anti-sparse coding formulations. As discussed above, if neocortex is organized in such a way to recursively solve the minimax optimization problem, recursive ANN search from low to high might end up with two equally possible solutions (x vs. $-x$).

Our final experiment deals with novelty detection in visual perception. In the so-called Thatcher effect [1], one can observe that it is relatively easy to detect inverted eyes and mouth from an upright face (see Fig. 4a). This detection of an anomaly (novelty) can be explained away because local patches around the eye and mouth regions are not compatible with the global impression of a human face. However, such novelty detection tasks become much more challenging with the inverted face (see Fig. 4b). Note that an inverted face is a rare event for HVS (e.g., face recognition becomes almost impossible [87]), and the hippocampus can only tell it is a novelty. Accordingly, top-down feedback in a hierarchical model will not detect the inconsistency between the inverted face and the eyes/mouth in normal position.



face image w/o and w. eyes inverted inverted face w/o and w. eyes inverted
 Fig. 4: Thatcher effect. It is much easier to detect inverted eyes from an upright face than from an upside-down face.

VI. TWO POSSIBLE IMPLEMENTATIONS USING ARTIFICIAL NEURAL NETWORKS

A. Asymmetrical Autoencoder: Decoder Matters More

The autoencoder represents a popular network architecture for unsupervised learning. A straightforward application of the sparse coding principle to the autoencoder is possible [69]. The hierarchical extension of the CSC inspires us to consider the design of an asymmetrical autoencoder, as shown in Fig. 5. Its asymmetrical design is also biologically inspired by H. Barlow's redundancy exploitation hypothesis [9]. The hierarchical organization of the neocortex is believed to reflect the nested structure of the physical world [35], indicating that the decoder (responsible for reconstruction of internal representations) plays a more important role than the encoder. More specifically, we can implement a prototype by combining

an over-parameterized autoencoder [80] with manifold-based novelty detection [77]. In the literature, autoencoder-based representation learning [88] is known to be capable of learning disentangled and hierarchical representations. Instead of storing training samples as attractors [80], we envision that a hierarchy of increasingly abstract concepts can become attractors of an over-parameterized asymmetrical autoencoder. In the presence of a new category (novelty), the escape mechanism will be activated and the decoder will be further expanded to accommodate the novelty class (i.e., consolidation of new memory).

Such unsupervised and continual learning can lead to monotonically increased memory capacity for associative memory implemented by overparameterized autoencoder. Unlike [80] treating sequence encoding as composition maps and limit cycles, we argue that a biologically more plausible mechanisms for memory storage is based on the rich interaction between sensory and motor systems. From the direct-fit perspective, motion dictates the regime within which the organism achieves the invariance to the composition of geometric transformations. The rich interaction between sensory and motor systems contributes to the formulation of reconstruction problems at multiple scales from parts to the whole [40] as an unsupervised mechanism of learning invariant representations. When the sensory motion goes out of the normal range (e.g., rotating a book continuously), the asymmetrical autoencoder will fail to recognize the object (i.e., it will be treated as a novelty).

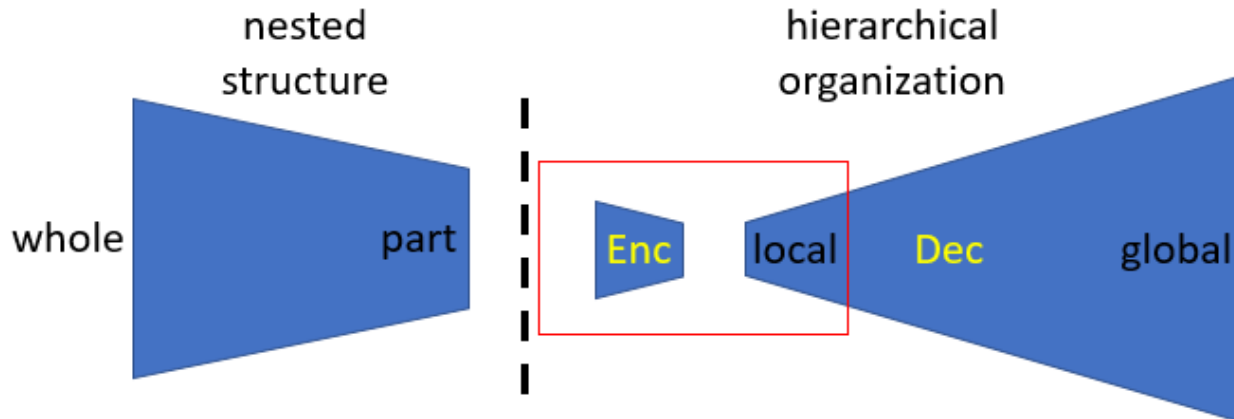


Fig. 5: Architecture of an asymmetrical autoencoder implementing the principle that nested structure in the physical world is reflected by the hierarchical organization of the neocortex (the dashed line marks the boundary between environment and organism; red box corresponds to the standard symmetric autoencoder).

B. Polychronous Neuronal Group (PNG) and Spiking Neural Networks

In [45], polychronization was conceived as a basic mechanism for computing with spikes. It is built upon Hebb's postulate, but extends it by relaxing the synchronous firing into polychronous time-locked patterns. Therefore, the group of neurons that are spontaneously organized by the fundamental process of spike-timing-dependent plasticity (STDP) is called Polychronous Neuronal Group (PNG). The mechanism of polychronization has recently been studied in [26, 44] as a plausible solution to the problem of feature binding. Using a spiking neural network (SNN), input training images (sensory stimuli) can be mapped to a hierarchy of PNGs by the emergence of polychronization. A mathematical abstraction of PNG is that it maps some partitioned space to a binary output (firing or no firing).

Unlike previous studies [26, 44], our over-parameterized direct fit model can be implemented on SNN with a divergent architecture, as shown in Fig. 6. Such a divergent architecture directly matches our intuition of abstracting complex cells by space composition. It is easy to see that the number of PNGs can grow exponentially as the number of neurons increases. Why do we need more PNGs at the higher level of the neocortex? The combinational coding argument (the end of Sect. III-B) suggests that PNGs could represent a biological plausible storage mechanism. They are physical implementations of attractors

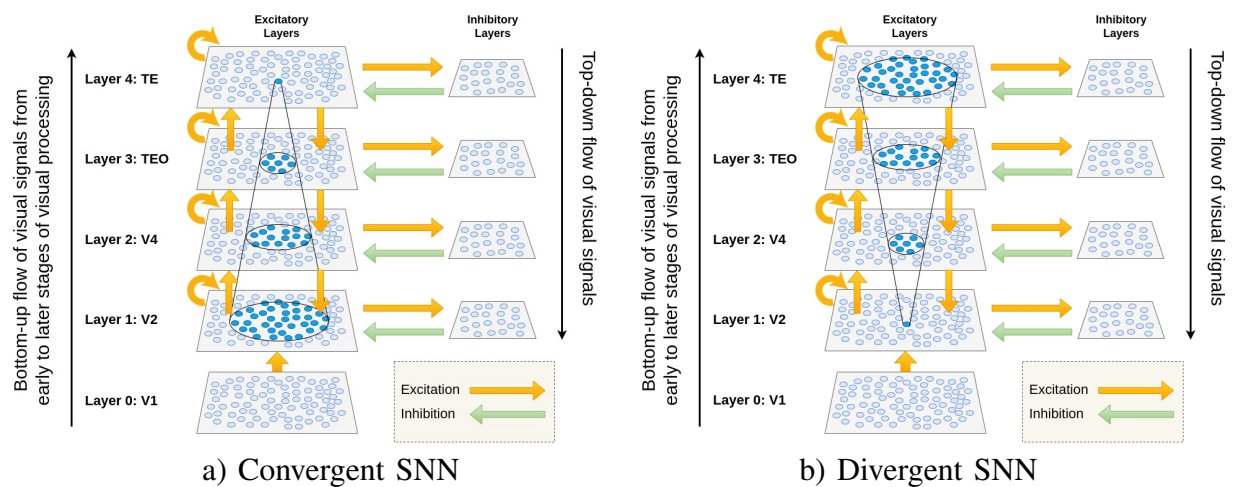


Fig. 6: Current implementation of SNN [26, 44] vs. the proposed implementation (there are more cells and synapses at higher levels than at lower levels).

in artificial neural networks. One way to experiment with our divergent SNN is to focus on its memory capacity; the other promising direction is to build a dorsal-ventral SNN to study the binding mechanism between “what” and “where” according to recent work of SNN for object vision [37] and spatial vision [75], respectively.

VII. CONCLUSIONS

We challenge the conventional way of modeling simple/complex cells by constructing space partitioning/composition in the topological space. Without imposing additional structures such as distance metric, we can construct product-manifold trees as a novel data structure suitable for the task of direct-fit visual perception. We demonstrate the biological plausibility of our construction and offer a sparse coding interpretation. It is possible to test the developed theory by constructing an asymmetrical autoencoder or a divergent SNN. For asymmetrical autoencoder, it will be interesting to study how a divergent architecture can maximize the capacity of associative memory. For divergent SNN, the objective is to experimentally verify the parallel and sequential binding mechanisms and their associated combination coding arguments, which might offer supporting evidence for the grandmother cell hypothesis.

REFERENCES

- [1] I. Adachi, D. P. Chou, and R. R. Hampton, “Thatcher effect in monkeys demonstrates conservation of face perception across primates,” *Current Biology*, vol. 19, no. 15, pp. 1270–1273, 2009.
- [2] E. H. Adelson, “Perceptual organization and the judgment of brightness,” *Science*, vol. 262, no. 5142, pp. 2042–2044, 1993.
- [3] E. H. Adelson and J. R. Bergen, “Spatiotemporal energy models for the perception of motion,” *Josa a*, vol. 2, no. 2, pp. 284–299, 1985.
- [4] C. C. Aggarwal, A. Hinneburg, and D. A. Keim, “On the surprising behavior of distance metrics in high dimensional space,” in *International conference on database theory*. Springer, 2001, pp. 420–434.
- [5] M. Aharon, M. Elad, and A. Bruckstein, “K-svd: An algorithm for designing overcomplete dictionaries for sparse representation,” *IEEE Transactions on signal processing*, vol. 54, no. 11, pp. 4311–4322, 2006.
- [6] Z. Allen-Zhu, Y. Li, and Z. Song, “A convergence theory for deep learning via over-parameterization,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 242–252.
- [7] S. Arora, N. Cohen, and E. Hazan, “On the optimization of deep networks: Implicit acceleration by overparameterization,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 244–253.
- [8] R. C. Atkinson and R. M. Shiffrin, “Human memory: A proposed system and its control processes,” in *Psychology of learning and motivation*. Elsevier, 1968, vol. 2, pp. 89–195.
- [9] H. Barlow, “Redundancy reduction revisited,” *Network: computation in neural systems*, vol. 12, no. 3, p. 241, 2001.
- [10] H. B. Barlow, “Unsupervised learning,” *Neural computation*, vol. 1, no. 3, pp. 295–311, 1989.
- [11] H. B. Barlow *et al.*, “Possible principles underlying the transformation of sensory messages,” *Sensory communication*, vol. 1, no. 01, 1961.

- [12] C. C. Bell, V. Z. Han, Y. Sugawara, and K. Grant, "Synaptic plasticity in a cerebellum-like structure depends on temporal order," *Nature*, vol. 387, no. 6630, pp. 278–281, 1997.
- [13] R. Bellman, "Dynamic programming," *Science*, vol. 153, no. 3731, pp. 34–37, 1966.
- [14] J. L. Bentley, "Multidimensional binary search trees used for associative searching," *Communications of the ACM*, vol. 18, no. 9, pp. 509–517, 1975.
- [15] J. S. Bowers, "On the biological plausibility of grandmother cells: implications for neural network theories in psychology and neuroscience," *Psychological review*, vol. 116, no. 1, p. 220, 2009.
- [16] N. Caporale, Y. Dan *et al.*, "Spike timing-dependent plasticity: a hebbian learning rule," *Annual review of neuroscience*, vol. 31, no. 1, pp. 25–46, 2008.
- [17] L. Chang and D. Y. Tsao, "The code for facial identity in the primate brain," *Cell*, vol. 169, no. 6, pp. 1013–1028, 2017.
- [18] D. Chen, X. Cao, F. Wen, and J. Sun, "Blessing of dimensionality: High-dimensional feature and its efficient compression for face verification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 3025–3032.
- [19] S. S. Chen, D. L. Donoho, and M. A. Saunders, "Atomic decomposition by basis pursuit," *SIAM review*, vol. 43, no. 1, pp. 129–159, 2001.
- [20] S. Chung and L. Abbott, "Neural population geometry: An approach for understanding biological and artificial neural networks," *Current opinion in neurobiology*, vol. 70, pp. 137–144, 2021.
- [21] J. Dai, H. Qi, Y. Xiong, Y. Li, G. Zhang, H. Hu, and Y. Wei, "Deformable convolutional networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 764–773.
- [22] S. Dasgupta and Y. Freund, "Random projection trees and low dimensional manifolds," in *Proceedings of the fortieth annual ACM symposium on Theory of computing*, 2008, pp. 537–546.
- [23] I. Daubechies, *Ten lectures on wavelets*. SIAM, 1992.
- [24] J. G. Daugman, "Complete discrete 2-d gabor transforms by neural networks for image analysis and compression," *IEEE Transactions on acoustics, speech, and signal processing*, vol. 36, no. 7, pp. 1169–1179, 1988.
- [25] D. L. Donoho, M. Elad, and V. N. Temlyakov, "Stable recovery of sparse overcomplete representations in the presence of noise," *IEEE Transactions on information theory*, vol. 52, no. 1, pp. 6–18, 2005.
- [26] A. Eguchi, J. B. Isbister, N. Ahmad, and S. Stringer, "The emergence of polychronization and feature binding in a spiking neural network model of the primate ventral visual system," *Psychological review*, vol. 125, no. 4, p. 545, 2018.
- [27] P. Földiak, "Forming sparse representations by local anti-hebbian learning," *Biological cybernetics*, vol. 64, no. 2, pp. 165–170, 1990.
- [28] D. François, V. Wertz, and M. Verleysen, "The concentration of fractional distances," *IEEE Transactions on Knowledge and Data Engineering*, vol. 19, no. 7, pp. 873–886, 2007.
- [29] J. H. Friedman, J. L. Bentley, and R. A. Finkel, "An algorithm for finding best matches in logarithmic expected time," *ACM Transactions on Mathematical Software (TOMS)*, vol. 3, no. 3, pp. 209–226, 1977.
- [30] J.-J. Fuchs, "Spread representations," in *2011 Conference Record of the Forty Fifth Asilomar Conference on Signals, Systems and Computers (ASILOMAR)*. IEEE, 2011, pp. 814–817.
- [31] K. Fukushima, "A hierarchical neural network model for associative memory," *Biological cybernetics*, vol. 50, no. 2, pp. 105–113, 1984.
- [32] J. J. Gibson, "The perception of the visual world." 1950.
- [33] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.
- [34] U. Hasson, S. A. Nastase, and A. Goldstein, "Direct fit to nature: an evolutionary perspective on biological and artificial neural networks," *Neuron*, vol. 105, no. 3, pp. 416–434, 2020.
- [35] J. Hawkins and S. Blakeslee, *On intelligence*. Macmillan, 2004.
- [36] S. Hayakawa and T. Suzuki, "On the minimax optimality and superiority of deep neural network learning over sparse parameter spaces," *Neural Networks*, vol. 123, pp. 343–361, 2020.
- [37] H. Hazan, D. Saunders, D. T. Sanghavi, H. Siegelmann, and R. Kozma, "Unsupervised learning with self-organizing spiking neural networks," in *2018 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2018, pp. 1–6.
- [38] D. Hebb, "The organization of behavior. a neuropsychological theory," 1949.
- [39] J. Hertz, A. Krogh, and R. G. Palmer, *Introduction to the theory of neural computation*. CRC Press, 2018.
- [40] G. Hinton, "How to represent part-whole hierarchies in a neural network," *arXiv preprint arXiv:2102.12627*, 2021.
- [41] G. Hinton and T. J. Sejnowski, *Unsupervised learning: foundations of neural computation*. MIT press, 1999.
- [42] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [43] D. H. Hubel, *Eye, brain, and vision*. Scientific American Library/Scientific American Books, 1995.
- [44] J. B. Isbister, A. Eguchi, N. Ahmad, J. M. Galeazzi, M. J. Buckley, and S. Stringer, "A new approach to solving the feature-binding problem in primate vision," *Interface focus*, vol. 8, no. 4, p. 20180021, 2018.
- [45] E. M. Izhikevich, "Polychronization: computation with spikes," *Neural computation*, vol. 18, no. 2, pp. 245–282, 2006.
- [46] H. Jégou, T. Furon, and J.-J. Fuchs, "Anti-sparse coding for approximate nearest neighbor search," in *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2012, pp. 2029–2032.
- [47] N. Kanwisher, J. McDermott, and M. M. Chun, "The fusiform face area: a module in human extrastriate cortex specialized for face perception," *Journal of neuroscience*, vol. 17, no. 11, pp. 4302–4311, 1997.
- [48] J. L. Kelley, *General topology*. Courier Dover Publications, 2017.
- [49] R. T. Knight, "Contribution of human hippocampal region to novelty detection," *Nature*, vol. 383, no. 6597, pp. 256–259, 1996.
- [50] Y. Le Cun and F. Fogelman-Soulié, "Modèle learning connectionists," vol. 2, no. 1, pp. 114–143, 1987.
- [51] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [52] M. Ledoux, *The concentration of measure phenomenon*. American Mathematical Soc., 2001, no. 89.
- [53] D. D. Lee and H. S. Seung, "Learning the parts of objects by non-negative matrix factorization," *Nature*, vol. 401, no. 6755, pp. 788–791, 1999.

- [54] J. M. Lee, B. Chow, S.-C. Chu, D. Glickenstein, C. Guenther, J. Isenberg, T. Ivey, D. Knopf, P. Lu, F. Luo *et al.*, “Manifolds and differential geometry,” *Topology*, vol. 643, p. 658, 2009.
- [55] D. A. Leopold and N. K. Logothetis, “Multistable phenomena: changing views in perception,” *Trends in cognitive sciences*, vol. 3, no. 7, pp. 254–264, 1999.
- [56] Y. Lyubarskii and R. Vershynin, “Uncertainty principles and vector quantization,” *IEEE Transactions on Information Theory*, vol. 56, no. 7, pp. 3491–3501, 2010.
- [57] S. G. Mallat, “A theory for multiresolution signal decomposition: the wavelet representation,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 11, no. 7, pp. 674–693, 1989.
- [58] C. v. d. Malsburg, “The correlation theory of brain function,” in *Models of neural networks*. Springer, 1994, pp. 95–119.
- [59] M. Markou and S. Singh, “Novelty detection: a review—part 1: statistical approaches,” *Signal processing*, vol. 83, no. 12, pp. 2481–2497, 2003.
- [60] —, “Novelty detection: a review—part 2:: neural network based approaches,” *Signal processing*, vol. 83, no. 12, pp. 2499–2521, 2003.
- [61] D. Marr, *Vision: A computational investigation into the human representation and processing of visual information*. MIT press, 2010.
- [62] D. Marr and W. T. Thach, “A theory of cerebellar cortex,” in *From the Retina to the Neocortex*. Springer, 1991, pp. 11–50.
- [63] D. Marr, D. Willshaw, and B. McNaughton, “Simple memory: a theory for archicortex,” in *From the Retina to the Neocortex*. Springer, 1991, pp. 59–128.
- [64] B. L. McNaughton, “Cortical hierarchies, sleep, and the extraction of knowledge from memory,” *Artificial Intelligence*, vol. 174, no. 2, pp. 205–214, 2010.
- [65] G. A. Miller, “The magical number seven, plus or minus two: Some limits on our capacity for processing information,” *Psychological review*, vol. 63, pp. 81–97, 1956.
- [66] P. M. Milner, “A model for visual shape recognition,” *Psychological review*, vol. 81, no. 6, p. 521, 1974.
- [67] V. Mountcastle, “An organizing principle for cerebral function: the unit module and the distributed system,” *The mindful brain*, 1978.
- [68] P. Nakkiran, G. Kaplun, Y. Bansal, T. Yang, B. Barak, and I. Sutskever, “Deep double descent: Where bigger models and more data hurt,” *Journal of Statistical Mechanics: Theory and Experiment*, vol. 2021, no. 12, p. 124003, 2021.
- [69] A. Ng *et al.*, “Sparse autoencoder,” *CS294A Lecture notes*, vol. 72, no. 2011, pp. 1–19, 2011.
- [70] B. A. Olshausen and D. J. Field, “Emergence of simple-cell receptive field properties by learning a sparse code for natural images,” *Nature*, vol. 381, no. 6583, pp. 607–609, 1996.
- [71] —, “Sparse coding with an overcomplete basis set: A strategy employed by v1?” *Vision research*, vol. 37, no. 23, pp. 3311–3325, 1997.
- [72] M. Page, “Connectionist modelling in psychology: A localist manifesto,” *Behavioral and Brain Sciences*, vol. 23, no. 4, pp. 443–467, 2000.
- [73] V. Pappas, Y. Romano, and M. Elad, “Convolutional neural networks analyzed via convolutional sparse coding,” *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 2887–2938, 2017.
- [74] V. Pappas, Y. Romano, J. Sulam, and M. Elad, “Theoretical foundations of deep learning via sparse representations: A multilayer sparse model and its connection to convolutional neural networks,” *IEEE Signal Processing Magazine*, vol. 35, no. 4, pp. 72–89, 2018.
- [75] F. Paredes-Vallés, K. Y. Scheper, and G. C. De Croon, “Unsupervised learning of a hierarchical spiking neural network for optical flow estimation: From events to global motion perception,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 42, no. 8, pp. 2051–2064, 2019.
- [76] J. Piaget and M. T. Cook, “The origins of intelligence in children.” 1952.
- [77] S. Pidhorskyi, R. Almoosen, and G. Doretto, “Generative probabilistic novelty detection with adversarial autoencoders,” *Advances in neural information processing systems*, vol. 31, 2018.
- [78] G. F. Poggio and L. J. Viernstein, “Time series analysis of impulse sequences of thalamic somatic sensory neurons,” *Journal of Neurophysiology*, vol. 27, no. 4, pp. 517–545, 1964.
- [79] R. Q. Quiroga, L. Reddy, G. Kreiman, C. Koch, and I. Fried, “Invariant visual representation by single neurons in the human brain,” *Nature*, vol. 435, no. 7045, pp. 1102–1107, 2005.
- [80] A. Radhakrishnan, M. Belkin, and C. Uhler, “Overparameterized neural networks implement associative memory,” *arXiv preprint arXiv:1909.12362*, 2019.
- [81] M. Riesenhuber and T. Poggio, “Hierarchical models of object recognition in cortex,” *Nature neuroscience*, vol. 2, no. 11, pp. 1019–1025, 1999.
- [82] F. Rosenblatt, “The perceptron: a probabilistic model for information storage and organization in the brain,” *Psychological review*, vol. 65, no. 6, p. 386, 1958.
- [83] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, “Learning internal representations by error propagation,” California Univ San Diego La Jolla Inst for Cognitive Science, Tech. Rep., 1985.
- [84] C. E. Shannon, “A mathematical theory of communication,” *The Bell system technical journal*, vol. 27, no. 3, pp. 379–423, 1948.
- [85] H. A. Simon, “The architecture of complexity,” *Proceedings of the American Philosophical Society*, vol. 106, no. 6, pp. 467–482, 1962.
- [86] P. Sterzer, A. Kleinschmidt, and G. Rees, “The neural bases of multistable perception,” *Trends in cognitive sciences*, vol. 13, no. 7, pp. 310–318, 2009.
- [87] J. W. Tanaka and M. J. Farah, “Parts and wholes in face recognition,” *The Quarterly journal of experimental psychology*, vol. 46, no. 2, pp. 225–245, 1993.
- [88] M. Tschannen, O. Bachem, and M. Lucic, “Recent advances in autoencoder-based representation learning,” *arXiv preprint arXiv:1812.05069*, 2018.
- [89] R. Vershynin, *High-dimensional probability: An introduction with applications in data science*. Cambridge university press, 2018, vol. 47.

- [90] R. Von der Heydt, E. Peterhans, and G. Baumgartner, "Illusory contours and cortical neuron responses," *Science*, vol. 224, no. 4654, pp. 1260–1262, 1984.
- [91] C. Von der Malsburg, "The what and why of binding: the modeler's perspective," *Neuron*, vol. 24, no. 1, pp. 95–104, 1999.
- [92] H. R. Wilson and J. D. Cowan, "Excitatory and inhibitory interactions in localized populations of model neurons," *Biophysical journal*, vol. 12, no. 1, pp. 1–24, 1972.
- [93] —, "A mathematical theory of the functional dynamics of cortical and thalamic nervous tissue," *Kybernetik*, vol. 13, no. 2, pp. 55–80, 1973.
- [94] J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, and Y. Ma, "Robust face recognition via sparse representation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 31, no. 2, pp. 210–227, 2008.
- [95] D. Yan, Y. Wang, J. Wang, H. Wang, and Z. Li, "K-nearest neighbor search by random projection forests," *IEEE Transactions on Big Data*, vol. 7, no. 1, pp. 147–157, 2019.
- [96] S. G. Young, K. Hugenberg, M. J. Bernstein, and D. F. Sacco, "Perception and motivation in face recognition: A critical review of theories of the cross-race effect," *Personality and Social Psychology Review*, vol. 16, no. 2, pp. 116–142, 2012.
- [97] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *arXiv preprint arXiv:1511.07122*, 2015.
- [98] F. Zenke, B. Poole, and S. Ganguli, "Continual learning through synaptic intelligence," in *International Conference on Machine Learning*. PMLR, 2017, pp. 3987–3995.
- [99] H. Zou, T. Hastie, and R. Tibshirani, "Sparse principal component analysis," *Journal of computational and graphical statistics*, vol. 15, no. 2, pp. 265–286, 2006.