# THE COOPERATIVE NETWORK ARCHITECTURE: LEARNING STRUCTURED NETWORKS AS REPRESENTATION OF SENSORY PATTERNS

**Pascal J. Sager**
Zurich University of Applied Sciences
University of Zurich
`sage@zhaw.ch`

**Jan M. Deriu**
Zurich University of Applied Sciences
`deri@zhaw.ch`

**Benjamin F. Grewe**
University of Zurich, ETH Zurich
`bgrewe@ethz.ch`

**Thilo Stadelmann**
Zurich University of Applied Sciences
European Centre for Living Technology (Venice, IT)
`stdm@zhaw.ch`

**Christoph von der Malsburg**
Frankfurt Institute for Advanced Studies
University of Zurich, ETH Zurich
`malsburg@fias.uni-frankfurt.de`

March 24, 2025

## ABSTRACT

We introduce the *Cooperative Network Architecture (CNA)*, a model that represents sensory signals using structured, recurrently connected networks of neurons, termed "nets." Nets are dynamically assembled from overlapping net fragments, which are learned based on statistical regularities in sensory input. This architecture offers robustness to noise, deformation, and out-of-distribution data, addressing challenges in current vision systems from a novel perspective. We demonstrate that net fragments can be learned without supervision and flexibly recombined to encode novel patterns, enabling figure completion and resilience to noise. Our findings establish CNA as a promising paradigm for developing neural representations that integrate local feature processing with global structure formation, providing a foundation for future research on invariant object recognition.

***Keywords*** net fragments · nets · neural code · pattern recognition · computer vision · neural networks · machine learning

## 1 Introduction

Artificial intelligence (AI) systems have demonstrated remarkable capabilities in processing diverse types of data, including images, text, and speech (LeCun et al., 2015; Sager et al., 2025). Yet, despite these advances, artificial vision systems continue to struggle with robust object representation and recognition, particularly under varying viewing conditions, occlusion, or noise (Qi et al., 2023; Simmler et al., 2021). This limitation highlights a fundamental gap between current AI systems and the flexible, context-aware perception observed in biological vision.

Existing approaches attempt to address these challenges through methods such as feature or sample engineering (Amirian et al., 2018; Moosavi-Dezfooli et al., 2016), domain adaptation (Csurka, 2017; Sager et al., 2022), and data augmentation (Carlucci et al., 2019; Tuggener et al., 2024). However, these techniques largely operate within the
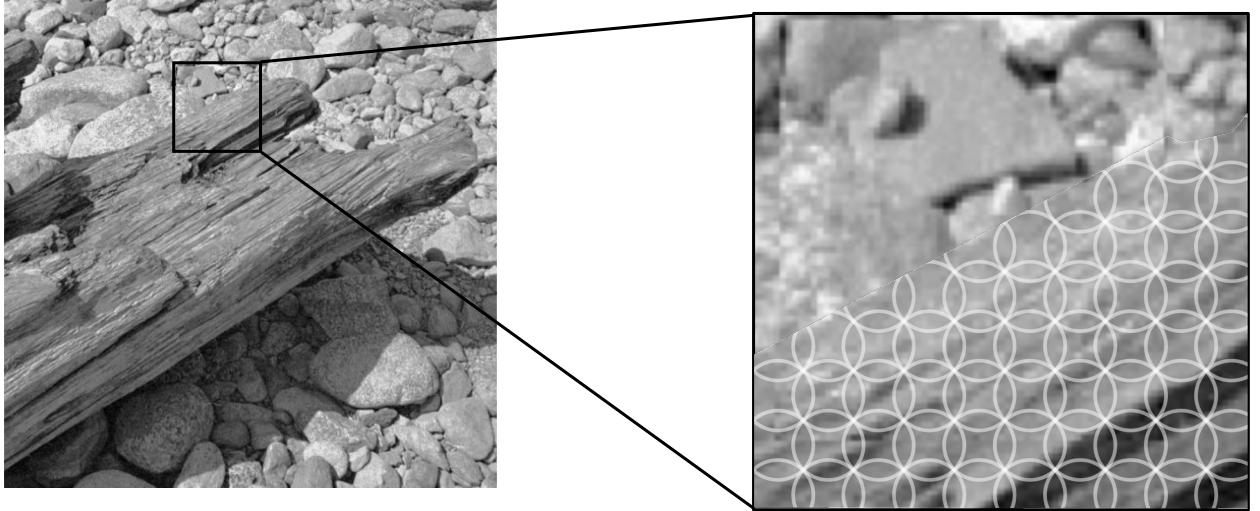
Figure 1: **Representation of Objects by Nets.** Objects are captured in their entirety by coherent nets composed of overlapping net fragments (schematically symbolized by white circles in the right panel). After training on natural images, each patch of visual space contains a complement of net fragments that represent textures that have been encountered with statistical significance. Net fragments overlap neuron-wise, and those activated by an object coalesce into a coherent net. The conundrum (stated in Olshausen and Field (2005)) that object contours can often not be found with the help of edge detectors due to lack of gray-level contrast (as in places inside the square in the left panel) may be resolved by the idea that contours are defined as borders of the nets covering objects (or covering the background). Figure adapted from Olshausen and Field (2005), with permission.

conventional paradigm of representing objects as static patterns of neuronal activations, often requiring large datasets and extensive supervision to achieve generalization.

In contrast, we propose a new paradigm by introducing the *Cooperative Network Architecture (CNA)*, a framework in which object representations emerge from dynamic interactions among recurrently connected neurons. Rather than encoding objects as distinct activation vectors, CNA organizes sensory information into "nets" – structured networks of mutually supporting neurons. These nets are assembled from overlapping "net fragments," which encode frequently co-occurring feature constellations and are acquired through Hebbian plasticity (Hebb, 1949). This compositional structure supports robust recognition, denoising, and figure completion through cooperative dynamics without external supervision.

In this paper, we present the first concrete implementation of CNA and demonstrate its application to a simple vision task. Our contributions include

- a formal mathematical definition of nets and net fragments,

- a computational framework for learning net fragments from sensory input,

- an illustration of how these fragments can be flexibly composed into object-spanning nets,

- and an analysis of the emergent filtering and generalization properties of nets.

Our results show that CNA offers a biologically inspired computational framework that integrates local and global information processing. It provides a foundation for more robust and adaptable perception systems, aligning with theoretical proposals for dynamic, network-based neural coding (von der Malsburg, 2018; von der Malsburg et al., 2022). By establishing the neural code and core computational mechanisms within a single representational area, this work provides the essential groundwork for future architectures in which multiple areas cooperate to support invariant object recognition – thus advancing toward the full potential envisioned in network-based neural coding models of perception.

## 2 The Net Fragment Framework

Before discussing the computational framework, we first introduce its building blocks: neurons, activations, nets, and net fragments. These definitions provide the basis for the discussions in the subsequent sections. For simplicity, we limit the definitions to a single area even though it could be extended to multiple inter-connected areas.

**Neurons.** A neuron $i$ is denoted by $n_i$. Each neuron $n_i$ has a *local neighborhood*, denoted by $N_r(n_i)$, which is the set of neurons $n_j$ such that the distance $d(n_i, n_j)$ between $n_i$ and $n_j$ is less than or equal to a threshold radius $r > 0$, i.e.,

$$N_r(n_i) = \{n_j \mid d(n_i, n_j) \leq r\}. \tag{1}$$

At each time step $t$, the *activity state* of a neuron $n_i$ is represented by $y_{i,t} \in \{0, 1\}$, where $y_{i,t} = 1$ indicates that the neuron is active at $t$. The activity of a neuron over time is characterized by an *activation pattern*, which is a sequence of activation states $\{y_{i,t}\}_{t=1}^{T}$, where $T$ is the total number of time steps (sequential processing steps).

**Connections.** Neurons interact via *connections*. A direct connection from neuron $n_i$ to another neuron $n_j$ is defined by a positive weight $w_{i,j} > 0$, which facilitates the propagation of neuronal activity from $n_i$ to $n_j$, so that the activity $y_{i,t}$ of neuron $n_i$ can influence the activity $y_{j,t+1}$ of neuron $n_j$ at the subsequent time step. In this work, we limit the connections to local neighborhoods:

$$\forall w_{i,j} \quad n_j \in N_r(n_i) \tag{2}$$

Additionally, we limit the weights to a range between 0 and 1, allowing neuron $n_i$ only to positively influence $n_j$:

$$w_{i,j} \in [0, 1] \tag{3}$$

We define a relation $n_i \sim n_j$ between two neurons to indicate that their activation patterns are highly correlated, i.e., $y_{i,t}$ often influences $y_{j,t+1}$. Under Hebbian learning (Hebb, 1949), the connection strength $w_{i,j}$ increases when there is $n_i \sim n_j$.

**Net Fragments.** A *net fragment* is defined as the set of neurons that provide direct activation support to a given neuron. In the context of this work, a neuron $n_j$ is said to provide *support* to another neuron $n_i$ if two conditions are met: (1) there exists a high correlation between two neurons $n_j$ and $n_i$ such that $n_j \sim n_i \Rightarrow w_{ji} > 0$; and (2) neuron $n_j$ was active at the previous time step, i.e., $y_{j,t-1} = 1$. Since synaptic weights are bounded within the interval $w_{j,i} \in [0, 1]$, the degree of support a neuron can provide is likewise restricted to this range.

Formally, the *net fragment* associated with neuron $n_i$ at time $t$ is given by

$$\mathcal{F}_t(n_i) = \{n_i\} \cup \{n_j \mid w_{j,i} > 0 \wedge y_{j,t-1} = 1\}. \tag{4}$$

In this formulation, any neuron $n_j \in \mathcal{F}(n_i)$ contributes to the activation likelihood of $n_i$ at time $t$. Notably, neuron $n_i$ is part of its own fragment, reflecting a form of self-support – if active at time $t$, the neuron may contribute to sustaining its own activation at the subsequent time step. Note that $n_i \sim n_j$ is not an equivalence relation since correlation is not transitive, meaning that one neuron can be part of multiple net fragments.

**Activation Function.** The decision of whether a neuron becomes active is based on the amount of support it receives from within the net fragment. We call this accumulated input the *membrane potential* at timestep $t$, which is a function $\mathcal{M}$ of the net fragment $\mathcal{F}_{t-1}(n_i)$ of a neuron $n_i$ at timestep $t - 1$.

$$a_{i,t}^{\mathcal{M}} = \mathcal{M}(\mathcal{F}_t(n_i)) \tag{5}$$

The membrane potential $a_{i,t}^{\mathcal{M}} \in \mathbb{R}$ is large when many other neurons provide support, meaning that they "confirm" the plausibility of a neuron's activity by being active themselves and exhibiting a connection that has been strengthened through previous mutual activity and Hebbian updates.

The activation of a neuron at timestep $t$ is a function of the membrane potential $a_{i,t}$ and environmental conditions $\mathcal{E}$:

$$y_{i,t} = \mathcal{A}(a_{i,t}, \mathcal{E}) \in \{0, 1\} \tag{6}$$

The environment conditions $\mathcal{E}$ include an *attenuation* signal, a constraint that selectively suppresses weak excitatory inputs by requiring neurons to accumulate a sufficiently large membrane potential before firing.

**Nets.** Neurons are not only influenced by direct connections within net fragments but also by indirect connections (e.g., if $n_i$ supports $n_j$ and $n_j$ supports $n_k$, then $n_i$ indirectly supports $n_k$). A set of directly and indirectly connected neurons is called a *net*, implementing the basic structure used in this work to represent objects. For simplicity in the following definitions, we introduce $\mathcal{F}_t^*(n_i)$, a net fragment $\mathcal{F}_t(n_i)$ in which $n_i$ is active:

$$\mathcal{F}_t^*(n_i) = \mathcal{F}_t(n_i) \quad \text{with} \quad y_{i,t} = 1 \tag{7}$$

To account for both direct and indirect influences on a neuron's activation, we use a recursive formulation of nets from the perspective of a single neuron. The base case is defined as the immediate net fragment:

$$\mathcal{N}_t^1[n_i] = \mathcal{F}_t^*(n_i) \tag{8}$$

For $r \geq 2$, the recursive relation is given by

$$\mathcal{N}_t^r[n_i] = \mathcal{N}_t^{r-1}[n_i] \bigcup_{n_j \in \mathcal{N}_t^{r-1}[n_i]} \mathcal{F}_t^*(n_j) \tag{9}$$

The complete network accessible from $n_i$ is then defined as the closure

$$\mathcal{N}_t^*[n_i] = \bigcup_{r=1}^{\infty} \mathcal{N}_t^r[n_i], \tag{10}$$

which comprises all neurons in the active subgraph that are reachable from $n_i$. This recursive formulation illustrates how local net fragments, which can be interpreted as elementary features, are progressively integrated to form higher-order features. Moreover, the explicit time dependence in these definitions underscores the dynamic nature of both the net fragments and the overall network, as they are continuously reassembled in response to the evolving support received from neighboring neurons, and increasing attenuation that prevents neurons with a low membrane potential from firing (see next section).

In the proposed architecture, no symmetry constraint is imposed on the weights, that is, $w_{i,j} \neq w_{j,i}$ in general. However, due to the temporal dynamics of the network, activations tend to stabilize over a longer time horizon $T$, such that for sufficiently large $t$, the activations become approximately stationary, i.e., $y_{i,t} \approx y_{i,t+1}$. Under these conditions, mutual influence between neurons becomes symmetric, meaning that $n_i \sim n_j$ implies $n_j \sim n_i$, and almost symmetric weights emerge naturally over time, even in the absence of explicit symmetry constraints.

As a result, it can be assumed that for any pair of distinct neurons $n_i \neq n_j$, their closures either converge to the same set of neurons or remain largely disjoint. Formally, this can be expressed as:

$$\mathcal{N}_t^*[n_i] \approx \mathcal{N}_t^*[n_j] \quad \text{or} \quad \mathcal{N}_t^*[n_i] \cap \mathcal{N}_t^*[n_j] \approx \emptyset. \tag{11}$$

## 3 Cooperative Network Architecture

In this section, we present a concrete instantiation of the proposed framework, specifically applied to the visual domain – although the framework is general and may be extended to other sensory modalities in the future. We refer to this implementation as the Cooperative Network Architecture (CNA). To realize our framework, we adopt a standard convolutional neural network (CNN) architecture (LeCun et al., 1989).

The CNA comprises two processing stages, depicted as yellow squares in Figure 2. The *stage 1 (S1)* performs fixed (i.e., not learned) feature extraction directly from the image input. The *stage 2 (S2)* implements the core principles of the described computational framework by learning net fragments.

The network dynamics unfold over discrete time steps $t = [0, \ldots, T]$. Since the feature extraction stage ($S1$) employs fixed filters, it yields a time-invariant activation pattern for a given input. In contrast, the second stage ($S2$) exhibits dynamic temporal evolution, where internal representations are iteratively refined at each time step ($t \rightarrow t + 1$). Specifically, we use a progressively increasing global attenuation signal that only allows neurons with a high membrane potential to remain active. During this dynamic processing, a neuron might be deactivated, whereby another neuron's membrane potential is decreased, triggering a chain reaction that can lead to the deactivation of incoherent nets.

While these refinements occur within the processing cycle of a single sample ($[0, \ldots, T]$), Hebbian learning adapts the weights (connections $w_{i,j}$ within net fragments) on longer timescales, wherein synaptic weight modifications occur between different sample-processing episodes, leading to structural plasticity in the network. The following sections discuss these two processing stages.
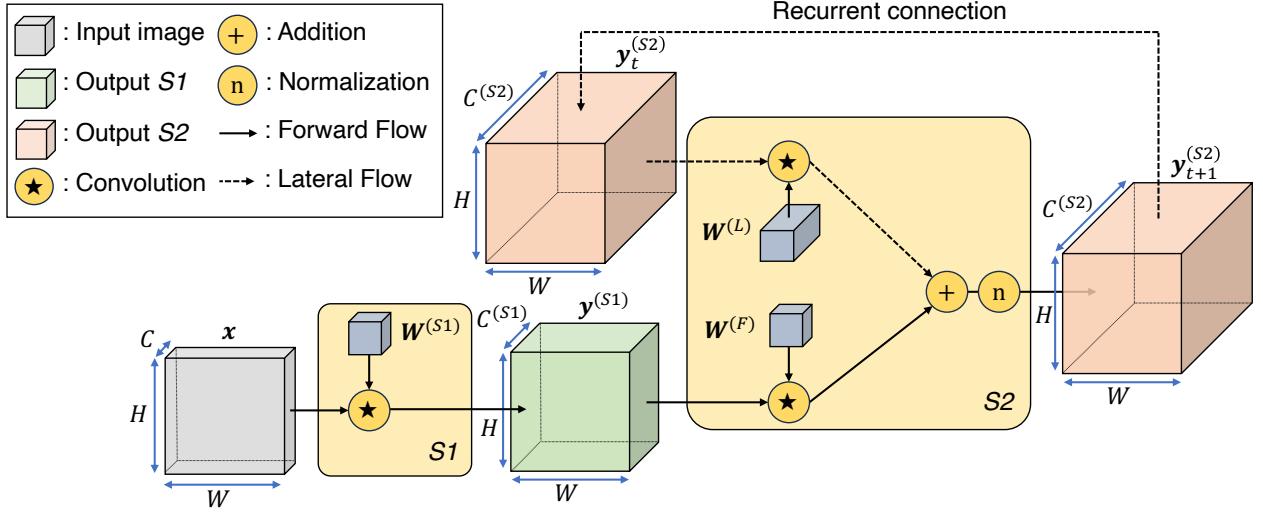
4

Figure 2: **The dynamics within the proposed CNA**: The input image is fed into stage $S1$ to obtain feature activations. The feature activations are then processed by stage $S2$, together with activations from neurons within the same layer.

### 3.1 Stage 1: Feature Extraction

Stage 1 is designed to model the sensory input by extracting initial neural activations from the raw image data. This extraction is essential due to the continuous nature of typical image inputs (e.g., grayscale or RGB), whereas the computational framework employed in this work operates on binary activations. Consequently, this stage serves as a preprocessing step, enabling the transformation of continuous-valued images into a suitable binary representation.

The input image is denoted as $\boldsymbol{x} \in \mathbb{R}^{C,H,W}$ (illustrated by the gray box in Figure 2), where $C$ indicates the number of image channels (1 for grayscale, 3 for RGB), while $H, W$ stand for the height and width of the image.

**Activation Computation.** To extract initial activation patterns, we use fixed, non-trainable filters (details see Appendix A). These filters are implemented as standard convolutional kernels, used to obtain the binary activation pattern $\boldsymbol{y}^{(S1)} \in \{0,1\}^{C^{(S1)},H,W}$ as follows:

$$\boldsymbol{y}^{(S1)}_{c_{\text{out}}} = \Theta\left(\sum_{c_{\text{in}}=0}^{C} \boldsymbol{W}^{(S1)}_{c_{\text{out}},c_{\text{in}}} \star \boldsymbol{x}_{c_{\text{in}}}\right) \tag{12}$$

Here, $\star$ denotes the convolution operator between the filters $\boldsymbol{W}^{(S1)} \in \mathbb{R}^{C^{(S1)},C,h,w}$ and the input image $\boldsymbol{x}$, while the parameters $h, w$ define the width and height of the convolutional filters, i.e., the spatial neighborhood from where the features are extracted. To binarize the activations, we use the heaviside function $\Theta$, which sets a neuron to one if its input is larger than zero, and to zero otherwise.

$$y_i = \Theta(a_i) = \begin{cases} 1 & \text{if } a_i > 0 \\ 0 & \text{otherwise} \end{cases} \tag{13}$$

Note that the width and height between input and output activations remain identical, but the number of channels changes from $C$ to $C^{(S1)}$, allowing the extraction of multiple features at each spatial location.

### 3.2 Stage 2: Net Fragments

$S2$ constructs net fragments from the binary feature activations $\boldsymbol{y}^{(S1)}$ of $S1$ and integrates them into coherent, spatially extended nets over time. To ensure stability, single neurons or small groups of neurons are suppressed through an attenuation mechanism, ensuring that only large-scale, consistent activation patterns persist.

5

**Network Topology.** Stage 2 uses convolutional layers to implement the connections $w_{j,i}$ within net fragments $n_j \in \mathcal{F}(n_i)$. Conceptually, this stage comprises a three-dimensional array of neurons with dimensions $C^{(S2)} \times H \times W$, where $C^{(S2)}$ denotes the number of feature channels in $S2$, and $H$ and $W$ correspond to the spatial dimensions of the input image.

The spatial distance between two neurons $n_i$ and $n_j$ is defined using the maximum norm (also known as Chebyshev distance), where adjacent neurons have a unit distance:

$$d(n_i, n_j) = \max\left(|h_i - h_j|, |w_i - w_j|\right) \tag{14}$$

Each neuron $n_i^{(S2)}$ in stage 2 is connected in two ways: (1) *Feedforward connections* to the corresponding neuron $n_i^{(S1)}$ in stage 1 and to its local neighborhood $Nr(n_i^{(S1)})$, and (2) *lateral connections* to neighboring neurons within stage 2 itself, denoted as:

$$Nr\left(n_i^{(S2)}\right) = Nr\left(n_i^{(S1)}\right) \cup \left\{n_j^{(S2)} \mid d(n_i^{(S2)}, n_j^{(S2)}) \leq r_2\right\} \tag{15}$$

where $r_2$ defines the radius of lateral connections. The convolutional architecture realizes these neighborhood structures via convolutional filters, where the filter size directly determines the radius of both feedforward and lateral connections. Using convolutional filters not only implements the connectivity restrictions to local neighborhoods but also allows detection of the same pattern at various positions, enabling learning of net fragments independent of position.

The forward connections from $S1$ to $S2$, denoted as $\boldsymbol{W}^{(F)}$, are necessary to activate cells within $S2$. It can be interpreted as a small group of neurons in $S1$ (based on the filter size of $\boldsymbol{W}^{(F)}$) supporting a single neuron within $S2$. In contrast, the lateral recurrent connections within $S2$, denoted as $\boldsymbol{W}^{(L)}$, are responsible for building the overlapping net fragments within $S2$, providing support amongst the neurons within $S2$. It is important to note that the lateral connection matrix $\boldsymbol{W}^{(L)}$ does not implement recurrent processing in the traditional sense of sequential token-wise updates, as found in RNN-based language models. Instead, $\boldsymbol{W}^{(L)}$ iterates over the same neuronal state vector $\boldsymbol{y}^{(S2)}$. Specifically, the lateral connections $\boldsymbol{W}^{(L)}$ calculate the state updates of $\boldsymbol{y}_{t \to t+1}^{(S2)}$ within an iterative refinement process in which individual neuronal activations $y_i^{(S2)} \in \boldsymbol{y}^{(S2)}$ may flip due to changes in support and attenuation.

**Membrane Potential.** To determine the support a single neuron $n_{i,t}$ receives from the net fragment $\mathcal{F}_{t-1}(n_{i,t-1})$, we first calculate the membrane potential $\mathcal{M}$ by applying convolutional operations on the forward and lateral connections and then taking the sum of these operations:

$$\boldsymbol{a}_{c_{\text{out}},t}^{\mathcal{M}(S2)} = \underbrace{\left(\sum_{c_{\text{in}}=0}^{C^{(S1)}} \boldsymbol{W}_{c_{\text{out}},c_{\text{in}}}^{(F)} \star \boldsymbol{y}_{c_{\text{in}}}^{(S1)}\right)}_{\textit{Forward Connection}} + \overbrace{\left(\sum_{c_{\text{in}}=0}^{C^{(S2)}} \boldsymbol{W}_{c_{\text{out}},c_{\text{in}}}^{(L)} \star \boldsymbol{y}_{c_{\text{in}},t-1}^{(S2)}\right)}^{\textit{Recurrent Connection}}. \tag{16}$$

Here, the forward connections have the shape $\boldsymbol{W}^{(F)} \in \mathbb{R}^{C^{(S2)},C^{(S1)},h,w}$ and the recurrent connections have the shape $\boldsymbol{W}^{(L)} \in \mathbb{R}^{C^{(S2)},C^{(S2)},h,w}$, where $C^{(S2)}$ is the number of channels in $S2$. Therefore, the membrane potential $\mathcal{M}$ at time $t$ is increased by the learned amount $w_{j,i}$ from all connected and active neurons ($y_{j,t-1} = 1$). Note that $\mathcal{F}_0(n_{i,0})$ only consists of neurons of $S1$ since all neurons in $S2$ are initialized to silent state ($y_{i,0} = 0$).

Higher support, indicated by larger $a_i^{\mathcal{M}(S2)}$, suggests a more likely pattern, as neurons encoding frequent patterns gain reinforcement not only from feedforward connections but also from lateral connections. Before we discuss how the weight matrices $\boldsymbol{W}^{(F)}, \boldsymbol{W}^{(L)}$ are initialized and trained, we describe how the membrane potential $\boldsymbol{a}^{\mathcal{M}(S2)}$ undergoes attenuation and is transformed into $\boldsymbol{y}^{(S2)}$.

**Attenuation.** The membrane potential is processed by an activation function $\mathcal{A}\left(\boldsymbol{a}^{\mathcal{M}}, \mathcal{E}\right)$. In our implementation, this function first normalizes the potential $\boldsymbol{a}^{\mathcal{M}} \to \boldsymbol{a}^{\text{norm}}$ to conform it into the range $[0, 1]$ (details see Appendix B), and then applies a globally increasing attenuation pressure $\boldsymbol{a}^{\text{norm}} \to \boldsymbol{a}$ to ensure that only neurons receiving substantial

support remain active. Specifically, attenuation is implemented as:

$$\boldsymbol{a}_{c,t}^{(S2)} = \left(\boldsymbol{a}_{c,t}^{\text{norm }(S2)}\right)^{\gamma}, \tag{17}$$

where $\gamma$ regulates the attenuation strength and increases over time according to $\gamma = \alpha + \beta t$ (when not stated differently, we set $\alpha = 1.2$ and $\beta = 0.2$), making activation progressively harder. This process leads to a sparsification of activations that occurs through both direct attenuation and reduced support from previously silenced neurons .

The final output is then binarized using the heaviside function $\Theta$:

$$\boldsymbol{y}_{c,t}^{(S2)} = \Theta\left(\boldsymbol{a}_{c,t}^{(S2)} - b^{(S2)}\right) \tag{18}$$

The activation at the last timestep $\boldsymbol{y}_{t=T}^{(S2)}$ represents a recurrently refined activation pattern where unsupported neurons have been suppressed, yielding a coherent representation.

**Alternative Neurons.** A single neuron that encodes a specific feature can be activated in multiple contexts. For example, a short straight-line segment may appear in various contexts – such as in different line drawings or digits. In our framework, each context is represented by a net, whereas individual features correspond to the activity of specific neurons.

To achieve *context disentanglement*, it is important to allocate sufficient representational capacity such that the same feature (e.g., a line segment) can be utilized across distinct contexts (e.g., different line drawings) without interference. This requires that a feature-representing neuron is capable of participating in multiple, context-specific net fragments. Without this capacity, a neuron would be forced to average over competing contextual associations, resulting in a non-specific or "blurry" representation.

Formally, a neuron representing a feature may be part of multiple net fragments: $n_i \in \mathcal{F}(n_j)$ and $n_i \in \mathcal{F}(n_k)$ where $n_j \neq n_k$). This implies that the neuron must form distinct lateral connections $w_{i,j}$ and $w_{i,k}$, which should not be concurrently active.

To address this requirement, we introduce "alternative neurons". Each alternative version of a neuron is activated by the same feature from $S1$ but has a distinct set of lateral connections in $S2$ to encode different contexts. This is implemented by increasing the channel dimension $C^{(S2)}$ of $\boldsymbol{W}^{(L)}$ by a factor $\kappa$, i.e.,

$$C^{(S2)} := C^{(S2)} \cdot \kappa \tag{19}$$

This expansion provides capacity for $\kappa$ alternative version of a neuron. Consequently, these alternative neurons receive identical bottom-up input from $S1$ while differing in the lateral support they receive in $S2$. A winner-take-all (WTA) competition mechanism ensures that each alternative neuron specializes in a distinct context, with only one being activated in any given instance.

The selection of the winning neuron is based on the correlation between feature patches and alternative kernel filters. Initially, patches of the same size as the feature kernel are extracted from each position in the input $\boldsymbol{y}_{t-1}^{(S2)}$. The correlation between these patches and the feature kernels is then computed, and at each location, the alternative kernel exhibiting the highest correlation is selected. This process ensures that the best-matching net fragment is activated, while all others remain suppressed.

At initialization, neuron selection is random. However, after Hebbian learning updates, the alternative neurons progressively specialize, leading to distinct and structured fragment representations. Over time, the most correlated alternative neuron is chosen for each input location, ensuring robust feature specialization across different contexts.

**Weight Initialization.** The forward connections $\boldsymbol{W}^{(F)}$ are initialized so that the activations of the $S1$ neurons are initially copied into all alternative $S2$ neurons in the same position. This is done by setting the connections at the center of the kernel (at position $(h/2, w/2)$) to 1 if $c_{in} = \lceil c_{out}/\kappa \rceil$, where $c_{in}$ denotes the index of the input channel, $c_{out}$ the index of the output channel, and $\kappa$ the number of alternative neurons. The recurrent connections $\boldsymbol{W}^{(L)}$ are initialized with zeros except for the self-coupling of neurons (recurrent connection to itself), which is set to 1.

**Learning of Weights.** Support within net fragments is *learned* via Hebbian plasticity (Hebb, 1949), so that frequently occurring patterns that lead to correlated activation patterns between neurons $n_i \sim n_j$, are "stored" in the system and receive more support:

$$\boldsymbol{W}^{(S2)} := \min\left(\max\left(\boldsymbol{W}^{(S2)} + \eta \cdot \boldsymbol{\rho}_{avg}, 0\right), 1\right) \tag{20}$$

7

Here, $\eta$ represents the learning rate and $\boldsymbol{W}^{(S2)}$ stands for both $\boldsymbol{W}^{(F)}$ and $\boldsymbol{W}^{(L)}$. Note that the same convolutional kernel is applied at various spatial positions. The decision to update this connection is based on the average correlation $\boldsymbol{\rho}_{avg}$ between all input and output neurons it connects (across all spatial positions). If the resulting average is positive, the corresponding connection increases as it connects more simultaneously active neurons than neurons that fire disjointly. Conversely, if the average correlation is negative, indicating more disjoint firing, the connection strength is reduced. Learning is applied after processing a sample (after $T$ timesteps) to improve stability.

## 4  Experiments

To demonstrate the compositionality of net fragments into nets, we train the model on binary images containing straight-line patterns so that the CNA learns net fragments corresponding to straight-line sub-patterns (smaller straight-line segments). However, for evaluation (after training and without weight updates), we use much more complex images such as line drawings of objects, and demonstrate that these complex patterns can be expressed as compositions (nets) of straight-line net fragments. This flexible compositional behavior contrasts strongly with conventional associative memory that cannot model such out-of-distribution data (see Section 6.1). To highlight the differences to ANNs, we train and evaluate an autoencoder on the same dataset and compare the representations.

We also examine robustness by introducing noise. Specifically, we assess the efficacy of filtering out random Gaussian noise and reconstructing partially occluded objects. Further experimental details are provided in the appendix: We describe the training and testing datasets in more detail in Appendix D, the feature extractor mechanism to obtain $\boldsymbol{y}^{(S1)}$ is described in Appendix A, the training details of our CNA model are in Appendix E, and the training details of the autoencoder in Appendix F.
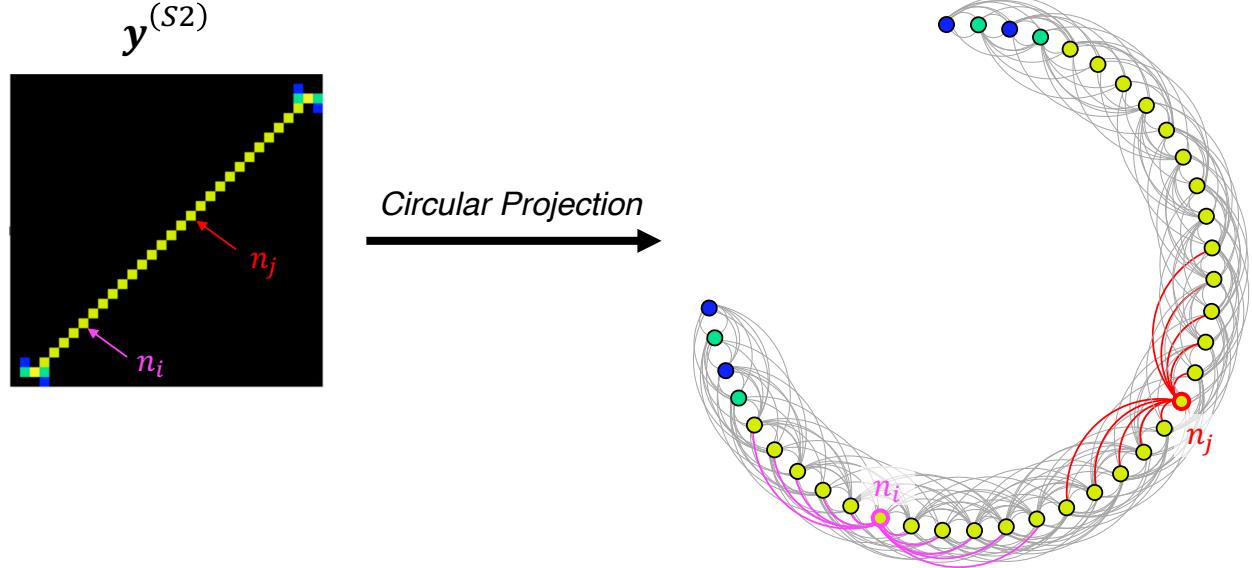
## 5  Results



Figure 3: The response of our model to a diagonal line input is displayed as a net of activated neurons. The left of the image shows the neuronal activation $\boldsymbol{y}^{(S2)}$ when observing the line, and the right visualizes the net formed by these neurons (activations $\boldsymbol{y}^{(S2)}$ with their connections $\boldsymbol{W}^{(L)}$). For convenience, the neurons are rearranged in a circle, and the net fragment (connections) of two randomly selected neurons, $n_i$ and $n_j$, are visualized in colors (purple and red).

Although the figures illustrate individual neuron activity, it is important to clarify that our objective is not to interpret patterns merely as isolated neuron activations, but rather as nets – cohesive, interconnected clusters of active neurons. To emphasize this perspective, Figure 3 illustrates the connectivity among activated neurons in $\boldsymbol{y}^{(S2)}$ (arranged in a circular layout for clarity), after feeding a diagonal line into the network. We select randomly two neurons $n_i$ and $n_j$ in the left panel and visualize their net fragments $\mathcal{F}(n_i)$ and $\mathcal{F}(n_j)$ based on $\boldsymbol{W}^{(L)}$ on the right. This figure highlights

that individual neurons (e.g., neurons $n_i$ and $n_j$) are embedded within broader nets of supportive connections. Without such network-based support, these neurons would be suppressed via attenuation mechanisms and would not appear in the left panel of Figure 3, which displays neuronal activity. Consequently, the subsequent visualizations should be understood as representing supported neurons embedded within nets, even though the underlying connectivity is not explicitly shown.
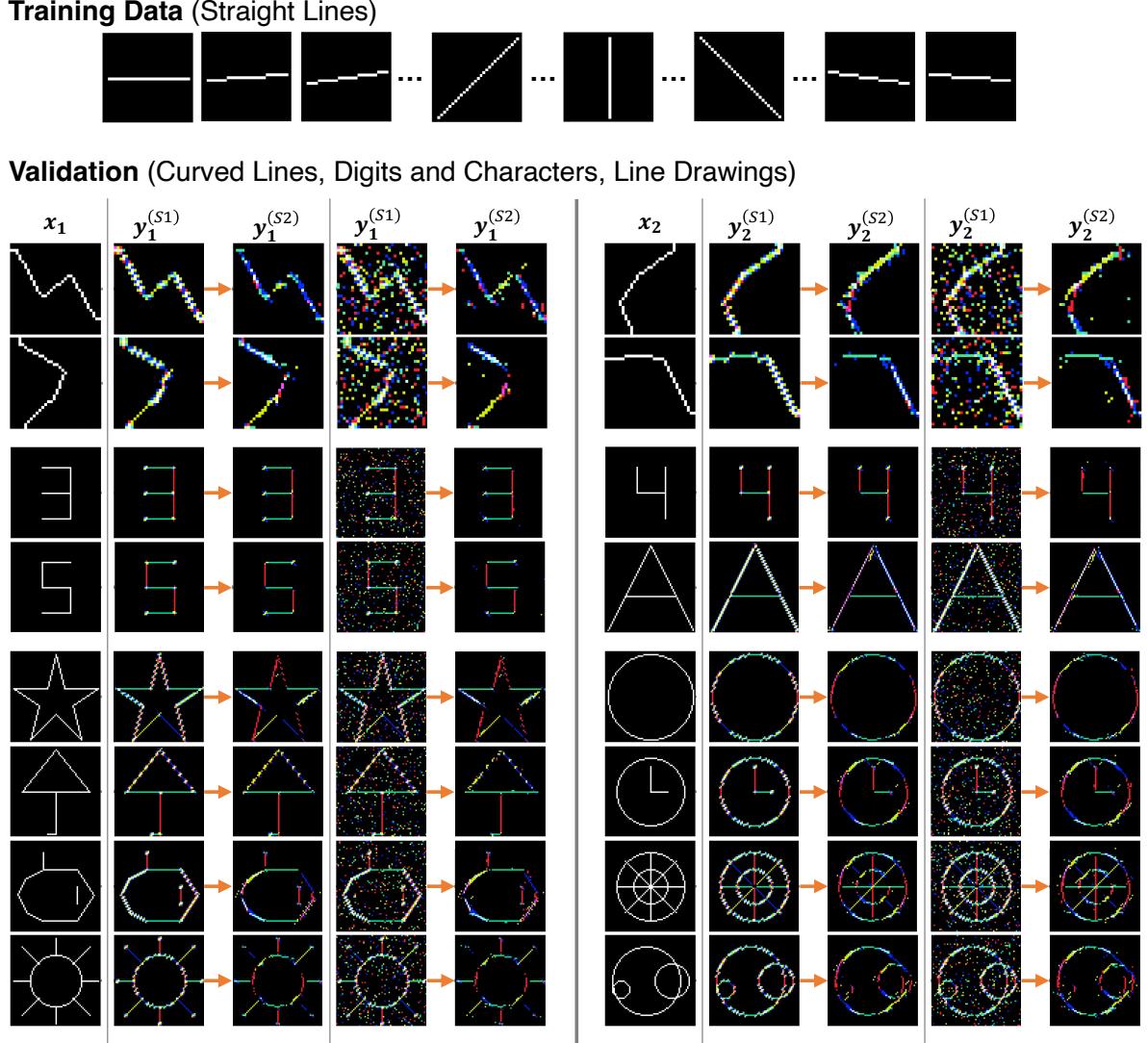
## 5.1 Compositionality



Figure 4: The top of the figure illustrates the straight lines used during the training phase. Below, various samples from the evaluation datasets are shown, including kinked lines (rows 1 and 2), digits and characters (rows 3 and 4), and line drawings (rows 5-8). For each row, two input examples, denoted as $x_1$ and $x_2$, are displayed alongside their respective feature activations, $y^{(S1)}$, and neuronal activations, $y^{(S2)}$, shown both without and with addition of Gaussian noise.

The model weights are trained using straight lines (see training data at the top of Figure 4). During inference – without additional learning – the learned net fragments are recursively composed to form nets representing significantly more complex structures (see validation data in rows 1-8 of Figure 4). Notably, although such structures are absent from the training set, the CNA successfully generalizes to represent kinked lines (rows 1 and 2), digits and characters (rows 3 and 4), and line drawings (rows 5–8).

9

These results demonstrate that while individual net fragments $\mathcal{F}_t$ must be learned (in our case we use Hebbian learning to train the connections within $\mathcal{F}_t$), the recursive composition described in eq. 9 to form $\mathcal{N}_t^*$ allows the integration of previously uncombined fragments. Consequently, the model exhibits combinatorial generalization through the structural recomposition of learned primitives.

The motivation for why we construct net fragments rather than just using the feature activation $\boldsymbol{y}^{(S1)}$ is demonstrated by introducing additional clutter: For each sample, we generate an input representation both with and without Gaussian noise. This noise corresponds to invalid patterns not observed during training (i.e., are not captured in net fragments), which lack sufficient recurrent support, leading to the deactivation of the corresponding neurons. Thus, the CNA model only allows valid sub-patterns to persist and be composed into novel structures, while invalid patterns, such as noise, are suppressed.
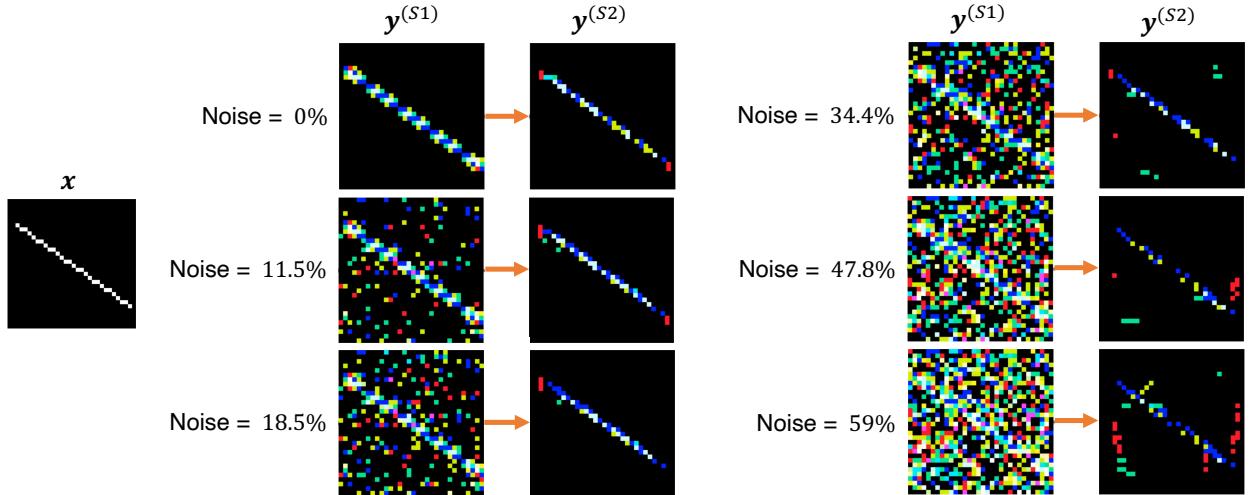
## 5.2 Filtering Gaussian Noise



Figure 5: Noise-corrupted feature activations $\boldsymbol{y}^{(S1)}$ alongside the corresponding output $\boldsymbol{y}^{(S2)}$, in which noise is reduced due to a lack of support within nets.

To further investigate noise filtering, we introduce Gaussian noise of varying degrees (from $0\%$ to $59\%$ noise at each spatial location - see Appendix D) to the afferent input in stage $S2$, thereby simulating an image structure that has already survived the filtering of feature kernels in $S1$ (thus increasing the difficulty, as the fragments must handle all the noise rather than relying on the feature extraction stage).

The results are visualized in Figure 5: While introducing up to $47.8\%$ noise, remarkably consistent outputs $\boldsymbol{y}^{(S2)}$ are generated. As additional noise is introduced, more undesired neurons receive support and persist in their activity. However, the overarching pattern remains distinguishable even when subjected to noise levels of up to $59\%$. In Appendix G, we quantitatively confirm the efficacy of noise filtering by measuring precision, recall, and noise filtration rate for various parameter settings, showing that lines remain well distinguishable for up to $59\%$ added noise.
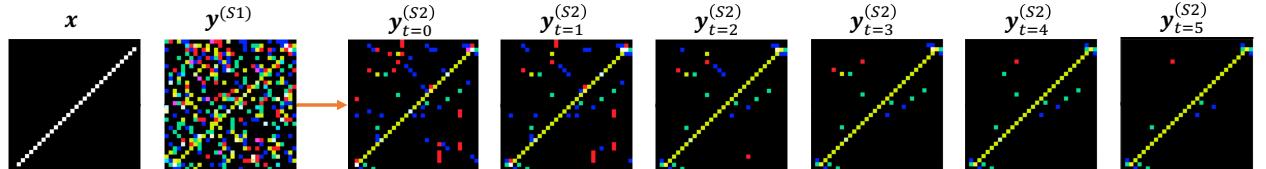


Figure 6: The filtering of (simulated) noise over time: Already at time step $t = 0$, most of the noise is suppressed. However, when attenuation increases over time, even more noise is filtered out. The plot is created after adding $10\%$ noise to each feature channel (using a bias $b^{(S2)} = 0.7$).

The reason for the filtering is that the additionally introduced clutter doesn't activate overlapping net fragments and cannot persist over time. Figure 6 depicts this systematic filtering across timesteps. The figure illustrates a substantial noise reduction already at time step $t = 0$ due to initial attenuation. In each subsequent time step, additional clutter is removed. This reduction is due to increased attenuation and falling support as a consequence of the drop-out of supporting neurons.
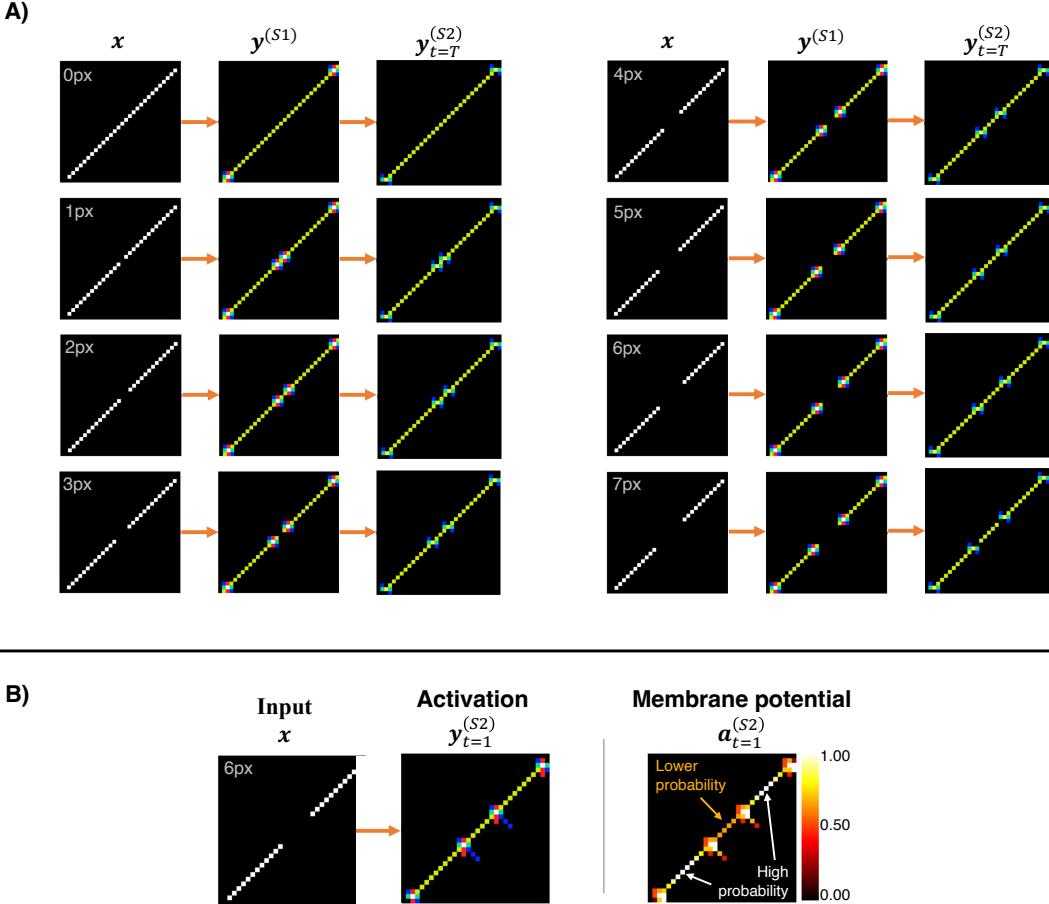
## 5.3 Reconstructing Figures



Figure 7: **Part A:** Input images of diagonal lines with different numbers of missing pixels, corresponding feature activations (middle column), and final output neuron activations. Activation bias $b^{(S2)} = 0.5$, attenuation coefficient $\gamma = 0.6 + 0.2t$. **Part B:** Diagonal line with 6 missing pixels in the center (shown on the left), the corresponding output activity (in the middle), and the membrane potential $a^{(S2)}$ displaying the "activation probabilities" (on the right) after the first time step.

If a fragmet $\mathcal{F}_t(n_i)$ contains a substantial number of (active) neurons, an inactive neuron $n_i$ can receive significant support, leading it to switch on and encouraging figure reconstruction ($\mathcal{F}_t(n_i) \rightarrow \mathcal{F}_t^*(n_i)$). Figure 7A shows how inactive neurons activate and thus demonstrates that nets can deal with occluded patterns. This reconstruction is reliable for up to 3 missing pixels (see Appendix G). In some cases, as in this Figure, the line is fully reconstructed for up to 6 removed pixels and partly reconstructed for 7 pixels.

**Membrane Potential Represents Expected Feature Activation.** In the presence of missing pattern elements, the membrane potential map $a^{(S2)}$ is higher at spatial locations where line features are observed and lower where they are not observed but expected (as shown in Figure 7B). Depending on the context, it may be desirable either to reconstruct the missing line segment – effectively bridging the gap between discontinuous features – or to preserve the separation

11

between distinct line elements. This behavior can be modulated by adjusting the hyperparameter $b^{(S2)}$ in eq. 18, which can govern the balance between feature completion and preservation of discontinuity, thereby allowing the model to selectively favor line reconstruction or separation.
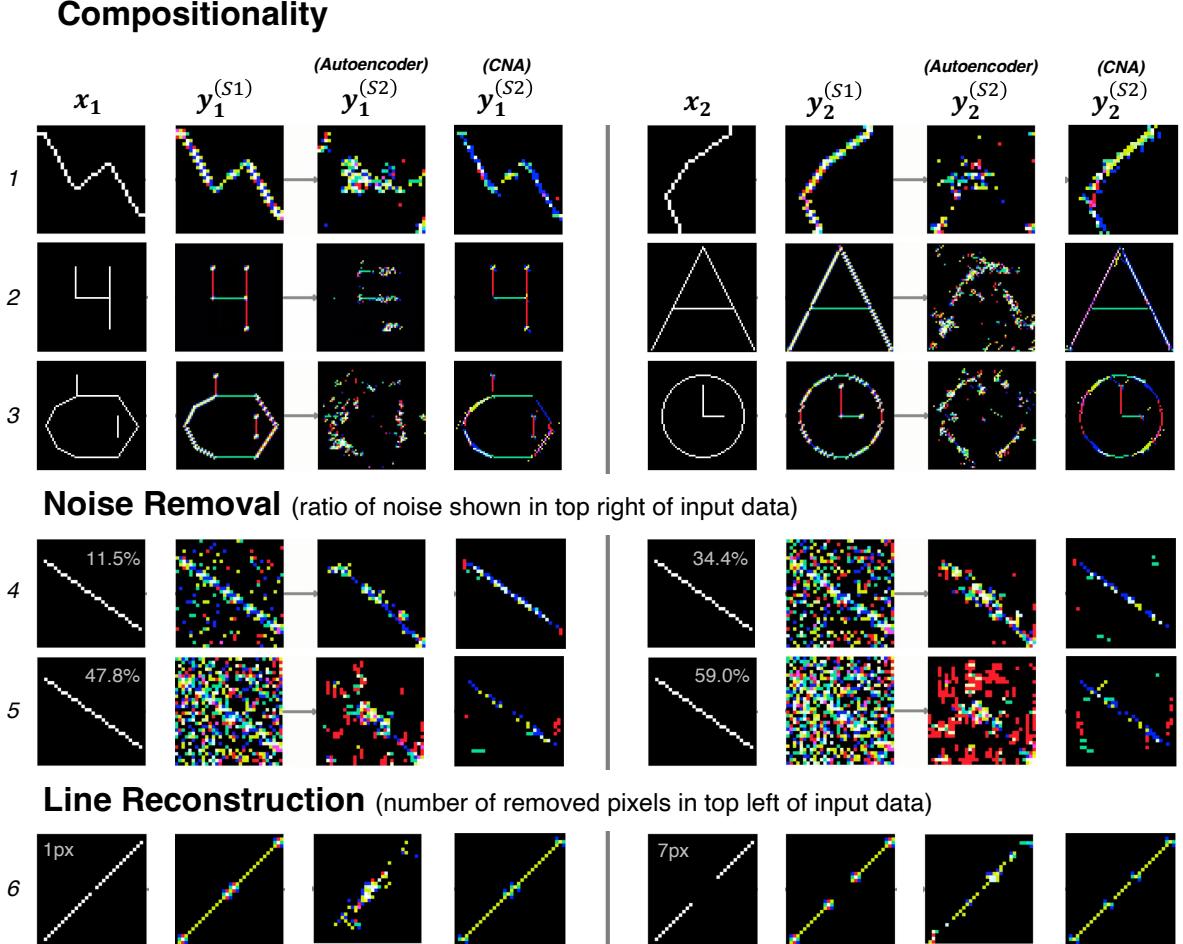
### 5.4 Comparison to an Autoencoder



Figure 8: Comparison of the outputs generated by the proposed CNA model and an autoencoder. Each row presents two samples, $x_1$ and $x_2$, alongside their respective feature activations and resulting representations (once for the CNA model and once for the autoencoder). Rows 1 to 3 illustrate how both models encode patterns that were not present in the training dataset. Rows 4 and 5 compare the models' performance in denoising tasks, while row 6 evaluates their ability to reconstruct lines.

Figure 8 compares performance between the proposed CNA model and an autoencoder (for model details, refer to Appendix F). While the autoencoder successfully learns to accurately represent line structures observed in the training data – achieving a mean square error (MSE) of less than $1 \times 10^{-3}$ between the input and the reconstruction by the decoder – it struggles with generalization to data not observed during training. For instance, rows 1 to 3 demonstrate that the CNA model is capable of representing previously unseen structures, while, in contrast, the autoencoder fails to capture such out-of-distribution patterns.

Rows 4 and 5 highlight the autoencoder's ability to represent line structures with up to 11.5% noise introduced. However, as the noise level increases, its performance decreases significantly. The CNA model, by contrast, remains robust under much higher noise levels, tolerating up to $59\%$ noise while still producing high-quality representations. Metrics on noise reduction, precision, and recall, presented in Appendix G, further substantiate the CNA model's superior performance.

In row 6, the models' capabilities for reconstructing missing line segments are compared. Both models can reconstruct the removed pixels, but the autoencoder introduces additional noise in regions other than the line center, where the pixels were removed. This observation is confirmed by the metrics in Appendix G: The autoencoder is slightly superior in reconstructing the removed pixels, but this comes at the cost of increased noise elsewhere, leading to lower overall precision and recall for figure reconstruction.

## 6  Relation to Alternative Concepts

To grasp the significance of the approach, it is important to compare it to alternative approaches.

### 6.1  Associative Memory

Associative memory (AM) (Amari, 1972; Hopfield, 1982) and the model presented in this work share the objective of retrieving previously encountered patterns, both utilize Hebbian synaptic plasticity for encoding information. However, fundamental differences exist in their structural and functional organization. In classical AM, neurons exhibit an all-to-all connectivity pattern, leading to the representation of structures as monolithic wholes, whereas our model constrains connectivity within learned net fragments. This restriction enables a more modular representation of complex patterns by dynamically assembling net fragments into structured nets, thereby allowing the formation of representations that extend beyond the training data distribution.

These differences arise from distinct learning mechanisms. In AM, a neural pattern is stored in a single event through the uniform strengthening of connections between co-activated neurons. Our approach, however, capitalizes on the observation that a large number of non-repeating patterns (as found in natural images) contain small recurring sub-patterns, such as texture patches. These recurring elements are captured within net fragments. The learning rule defined in eq. 20 selectively reinforces connections corresponding to the most relevant sub-patterns. Specifically, frequent sub-patterns exhibit a consistently positive average correlation $\rho$, which leads to the progressive increase of their associated synaptic weights. In contrast, rare or non-informative patterns exhibit low or even negative average correlations, causing their associated weights to decay toward zero. This Hebbian mechanism effectively filters out noise and ensures that only the most salient and statistically dominant sub-patterns are preserved in the learned representation.

A fundamental limitation of AM is the interference between stored patterns unless they are statistically independent (i.e., orthogonal), meaning that minimal neuron-wise overlap exists. Sparse coding, as proposed by Tsodyks and Feigel'man (2007) or Palm (2013), mitigates interference by activating only a small subset of neurons. Our model offers a more intrinsic solution by representing each input feature using multiple alternative neurons. These neurons possess identical receptive fields (as observed by Cossell et al. (2015)), but can exhibit different excitatory outputs. This redundancy enables overlapping net fragments to reduce neuron-wise interference by utilizing distinct alternative neurons.

### 6.2  Boltzmann Machine

Key aspects of our model become clearer when compared to the Boltzmann Machine (BM) introduced by Hinton and Sejnowski (1983). Both models allow recurrent connections (as also proposed in Mumford (2002)), but they also exhibit significant differences. Nevertheless, as our focus is limited to one area with net fragments ($S2$) and is not the full object detection architecture as proposed by von der Malsburg et al. (2022), we highlight only select differences here.

The BM makes a clear distinction between visible and hidden ("latent") neurons. In the CNA this distinction is softened: it posits neurons that are activated by sensory input, but of which the majority are quickly silenced by attenuation in favor of those neurons that support each other through previously learned connections. A winning alternative neuron can be interpreted as the hypothesis not only that the neuron has detected its trigger pattern in its afferents but also that the neuron has connections to a number of other active neurons and is consistent with them in light of previous experience.

In both BM and our model, the activity of hidden, or in our case alternative, neurons, is controlled by the same type of dynamic equations, with the difference being that we do not insist on symmetric connections and therefore do not derive an energy function (see, however, the remark made in Hinton and Sejnowski (1983), that given sufficiently redundant support of active neurons, the absence of up to half of the connections due to lack of symmetry does not jeopardize the stability of a state). However, even without an energy function, convergence is not an issue, as alternative neurons undergo winner-take-all competition.

### 6.3 Artificial Neural Networks (ANN)

Artificial neural networks (ANNs) rely on various architectures and learning techniques, including multilayer perceptrons (MLPs) (Prince, 2023), transformers (Vaswani et al., 2017), and autoencoders (Hinton and Salakhutdinov, 2006). These systems process sensory input through multiple layers to generate a neuronal representation. A defining characteristic of ANN systems is their strict feedforward architecture, where connections from input to representation and from representation to output are trained via backpropagation of error (Rosenblatt, 1962, p. 292).

In contrast, the neuronal representation within the CNA consists of nets that activate as a whole, with individual neurons depending on the activation of others in the same net. Thus, while ANNs rely on a global structuring of neuronal activations, the CNA relies on a local self-organizing process.

Moreover, ANNs rely on hierarchical layers where neurons respond to features of increasing complexity. The CNA replaces this hierarchy with a single layer, where net fragments of different sizes capture features of varying complexity. These fragments interconnect to form larger nets that represent entire objects, naturally solving the binding problem, whereas ANNs do not encode explicit relationships between features, leading to ambiguities such as hallucinated feature conjunctions.

Furthermore, input representation in the CNA emerges from statistical regularities in the input, making it independent of specific tasks. In contrast, ANNs bias feature extraction to serve predefined training objectives, which limits their adaptability when retraining for different functions.

## 7 Discussion and Outlook

The CNA presented here establishes a novel approach for constructing coherent nets as representations of input data. This framework enables the encoding of complex patterns through the composability of previously learned sub-patterns (net fragments). Compared to existing neural coding frameworks, CNA offers distinct advantages, particularly in terms of robustness – demonstrated through its ability to filter noise and reconstruct figures – and its capacity to generalize to out-of-domain data, as evidenced by its capability to represent previously unseen objects.

For future research, we outline two key directions: First, an essential next step is the extension of CNA to natural image domains. This would allow net fragments to capture and encode texture information encountered in real-world scenarios, thus realizing the ambition expressed in Figure 1. Second, the current implementation of CNA is restricted to a single processing area with net fragments. Following the theoretical framework proposed by von der Malsburg et al. (2022), the architecture could be expanded to incorporate multiple areas, thereby enabling object recognition while achieving invariance to position, scale, and in-plane rotation. Object recognition could be realized using structure-sensitive, homeomorphic mappings between nets (von der Malsburg, 1981; Anderson and van Essen, 1987; Olshausen et al., 1995; Wiskott et al., 1997; Wolfrum et al., 2008; Memisevic and Hinton, 2010; Arathorn et al., 2013), which would facilitate the mapping of input representations onto a more stable schematic representation of an object class or an exemplar of the object to be recognized. The adoption of homeomorphic mappings, as opposed to alternative methods, preserves the core principle of nets and net fragments within and across multiple processing areas, thereby maintaining structural coherence in representation.

The proposed neural coding framework not only advances representational robustness but also aligns with principles of self-organization observed in cortical processing. Its ability to learn and flexibly recombine net fragments without explicit supervision suggests a pathway toward more adaptive and interpretable neural representations. Future extensions to multi-area architectures could bridge the gap between bottom-up feature learning and top-down predictive coding. These developments position CNA as a foundational step toward more structured and generalizable neural computation.

## Funding

## Data Availability Statement

The code to generate the data as well as to reproduce the results can be found after the publication of this work on Github at `https://github.com/sagerpascal/cna`.
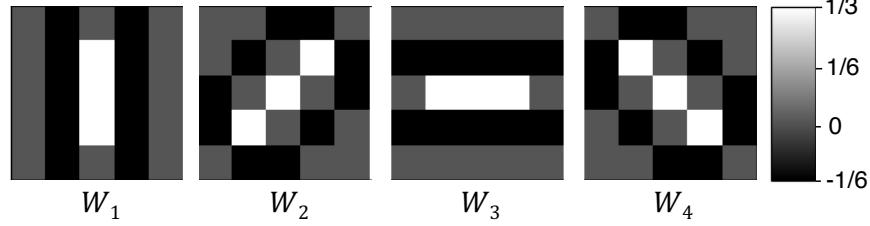
## Appendix

## A  Feature Extraction



Figure 9: The hand-crafted convolutional kernels of the feature extractor.

The stage $S1$ serves to extract initial binary features from the images. The input images have $C = 1$ channel and $C^{(S1)} = 4$ features are extracted at each spatial location. We utilize hand-crafted filters with a size of $5 \times 5$ as displayed in Figure 9, each of the four filters corresponding to a particular line orientation (vertical, positive diagonal, horizontal, and negative diagonal). This approach is motivated by the observation that straight lines when viewed locally, can be expressed as combinations of these fundamental line types (e.g., a line with a $20°$ angle activates the horizontal and the $+45°$ diagonal feature neurons). The use of such manually designed filters allows a more straightforward interpretation of the results. While more complex datasets necessitate learned filters, our primary focus in this work is on constructing net fragments, and thus, we favor the simplicity of hand-crafted filters.

These filters move across the image with a step size ('stride') of 1 (assuming 0 as input signal beyond the image border: 'zero padding') so that the input image and output image are of the same size. Thus, the input $x$ in the feature extraction stage are images of size $(1 \times 32 \times 32)$ and the output $y^{(S1)}$ feature activation maps of size $(4 \times 32 \times 32)$. The dataset's lines activate multiple filters at different positions, facilitating the construction of net fragments across these feature channels. We convert the floating-point output of the convolutional weights to binary activations using the heaviside function $\Theta$ with a bias of $b^{(S1)} = 0.5$, i.e. neurons fire a binary spike if $a_j > 0.5$ (see eq. 12).

## B  Normalization

The membrane potential in $S2$ is normalized using a two-step approach that confines the neuronal activations $a^{\mathcal{M}(S2)}$ to the interval $(0, 1)$.

We first apply a size limit:

$$a_{c,t}^{\mathcal{M}(S2)} = \begin{cases} a_{c,t}^{\mathcal{M}(S2)}, & \text{if } a_{c,t}^{\mathcal{M}(S2)} \leq \lambda \\ \lambda - \frac{1}{2}(a_{c,t}^{\mathcal{M}(S2)} - \lambda), & \text{otherwise.} \end{cases} \tag{21}$$

According to this formula, $a^{\mathcal{M}}$ grows undiminished with $a^{\mathcal{M}}$ until it reaches the value $\lambda$, and thereafter actually diminishes with slope $-1/2$. (In our experiments we found the value $\lambda = 1.3 \cdot \frac{h+w}{2}$ to work well). This saturation helps to mitigate imbalances between net fragments with differing numbers of participating neurons. In the next step, we introduce the normalization

$$a_{c,i,t}^{\text{norm}(S2)} = \max \left( \frac{a_{c,i,t}^{\mathcal{M}(S2)}}{\max_{i'} a_{c,i',t}^{\mathcal{M}(S2)}}, 0 \right) \tag{22}$$

Here, $i$ refers to spatial locations within the image, i.e., $i = (h, w)$ where $h \in H, w \in W$. This normalization confines activations to the interval $[0, 1]$, where neurons receiving maximal support are mapped to 1, while those with less support are mapped to correspondingly smaller values. This normalization is relative to the maximal value that itself is growing during learning as it is important to give neurons initially (when recurrent connections are still small) a chance to fire and later when recurrent weights between correlated patterns have been learned, to suppress those neurons that are part of noise patterns.

## C   Single Weight-Kernel Implementation

So far, we introduced the weights $\boldsymbol{W}^{(F)}$ and $\boldsymbol{W}^{(L)}$ in $S2$ as two independent convolutional kernels. For simplicity in our algorithm, and since we use the same kernel size for the forward and recurrent connections, we stack the two kernels as $\boldsymbol{W}^{(S2)} = (\boldsymbol{W}^{(F)}; \boldsymbol{W}^{(L)})$ as well as the inputs $\boldsymbol{y}^{(S1)}$ and $\boldsymbol{y}_{t-1}^{(S2)}$, leading to the equivalent formulation of eq. 16:

$$\boldsymbol{a}_{c_{\text{out}},t}^{\mathcal{M}(S2)} = \sum_{c_{\text{in}}=0}^{C^{(S2)}} \boldsymbol{W}_{c_{\text{out}},c_{\text{in}}}^{(S2)} \star \left( \boldsymbol{y}_{c_{\text{in}}}^{(S1)}; \boldsymbol{y}_{c_{\text{in}},t-1}^{(S2)} \right) \tag{23}$$

## D   Dataset

The dataset used for training comprises binary images measuring ($32 \times 32$) pixels, each sample depicting a straight line going through the image center, starting and ending 2 pixels from the image boundary (leading to 59 distinct lines of different angles). The images are generated when required and may vary from one epoch to another. During each training cycle, we randomly generate 300 image instances (each line about 5 times). During evaluation, we sample each of the 59 distinct lines exactly once. In contrast to training, the samples used for evaluation comprise local distortions in the form of additive Gaussian nose and missing line segments (subtractive noise), and we evaluate the net fragments' capability to remove these.

To introduce additive Gaussian noise to the afferent input in stage $S2$, we probabilistically flip neurons within each channel with a probability of up to 20%. Since we use four feature channels in our experiments, this corresponds to a probability of $1 - (1 - 0.2)^4 = 59\%$ that a neuron is flipped at any spatial location. To evaluate subtractive noise, we create discontinuous lines by deleting a line segment of 1 to 7 pixels in the middle.

To demonstrate the generalization capability of CNA, we generate validation samples that differ from the ones used during training. Specifically, we produce binary images measuring ($32 \times 32$) pixels depicting kinked lines by generating four random points within the image and drawing straight lines between these points.

Additionally, we generate binary images measuring ($64 \times 64$) pixels by hand that represent more complex structures. These images include line drawings of digits 0-9, characters, and objects such as a house, a star, a wheel, a clock, etc. Some of these images are shown in Figure 4, and an extensive list can be found at `https://github.com/sagerpascal/cna`.

## E   Parameters

We train the model with a learning rate of $\eta = 0.2$ for 100 training cycles (epochs). The model uses $\kappa = 10$ alternative neurons with recurrent connections spanning a kernel of $11 \times 11$ neurons. Since we use $C^{(S1)} = 4$ filters for feature extraction in $S1$ (see Appendix A), the number of alternative channels in $S2$ corresponds to $C^{(S2)} = 4 \cdot \kappa = 40$. The resulting weight matrices have dimensions $\boldsymbol{W}^{(F)} \in \mathbb{R}^{40 \times 4 \times 11 \times 11}$ and $\boldsymbol{W}^{(L)} \in \mathbb{R}^{40 \times 40 \times 11 \times 11}$, respectively $\boldsymbol{W}^{(S2)} \in \mathbb{R}^{40 \times 44 \times 11 \times 11}$.

## F   Autoencoder Parameter

We compare our CNA model to an ANN, specifically an autoencoder (Hinton and Salakhutdinov, 2006), which is a widely used self-supervised technique for generating input representations. An autoencoder is particularly suited for comparison because it includes a built-in decoder, allowing for feature reconstruction and visualization, which is required for feature comparison. To ensure a fair comparison, we also employ convolutional filters (LeCun et al., 1989) to the autoencoder and incorporate state-of-the-art design principles (we even permit more data and feature channels to the autoencoder).

The autoencoder is trained to reconstruct the feature activations, $\boldsymbol{y}^{(S1)}$, using the same training dataset as our CNA model. It consists of an 8-layer architecture, with an encoder composed of convolutional layers with a kernel size of 3 and a stride of 2, with progressively increasing channel counts of 32, 64, 128, and 256 channels. The decoder mirrors this structure, using transposed convolutional layers with 256, 128, 64, and 32 channels, respectively. ReLU activations are applied after each layer. Consequently, the autoencoder is a relatively large model with 777k parameters, utilizing up to 256 channels, whereas the CNA model operates with only 44 channels.

16

The autoencoder is trained by minimizing the L2 norm between the input features and their reconstructed counterparts, using the Adam optimizer (Kingma and Ba, 2015) with a learning rate of $1 \times 10^{-4}$, beta parameters of $(0.9, 0.999)$, and a weight decay of $1 \times 10^{-8}$.

To provide a robust comparison, we increase the number of training samples for the autoencoder. Rather than training for 100 epochs with 300 samples per epoch as done for the CNA model, we train the autoencoder for 200 epochs, using $10{,}000$ images per epoch. This results in a significantly larger training set but ensures that the autoencoder achieves a highly accurate reconstruction. The autoencoder is trained with a batch size of $512$.

## G   Quantifying Noise Filtration

Our evaluations focus on the efficacy of net fragments in suppressing added Gaussian noise and complementing missing or occluded patterns. For both evaluation experiments, we report *recall* and *precision*: Recall is defined as the proportion of neurons that maintain activity (on) when noise is added, and precision quantifies the extent to which neurons activated in the presence of noise were also active without noise. To assess the system's noise reduction capability in the case of Gaussian noise, we additionally measure the percentage of flipped neurons that revert to their original activation state, referring to this measurement as the *noise reduction rate*. To assess subtractive noise, we measure the similarity between the original (created without noise) and reconstructed (created with noise) feature activations at the spatial locations where pixels have been removed. We call this metric the *feature reconstruction rate*.

**Filtering Gaussian Noise.**   Figure 10 displays noise reduction rates, recall, and precision across different parameter configurations for different levels of added noise. The metrics demonstrate that increasing the attenuation coefficient or the activation bias results in enhanced noise filtration.

The overall efficacy of noise reduction remains consistently high across most parameter configurations, exhibiting a constant noise reduction rate of $\geq \approx 95\%$. This rate persists regardless of the magnitude of the introduced noise. As the noise suppression rate is approximately constant, this means that if there is more noise in the input, there will also be more noise in the output, albeit limited to around 5% of the input noise.

Furthermore, all CNA models (except the one with the lowest power factor and bias) outperform the autoencoder, demonstrating that they achieve performance that is competitive with state-of-the-art deep learning models.

**Reconstructing Subtractive Noise.**   Figure 11 displays the feature reconstruction rate, precision, and recall for removed line fragments of different lengths. Having a lower activation bias or attenuation coefficient $\gamma$ leads to a better reconstruction of removed patterns, as these settings allow neurons with less support to activate.

Compared to the autoencoder baseline, the CNA model achieves higher recall and precision but a lower feature reconstruction rate. This indicates that the autoencoder is superior at reconstructing the pixels that have been removed but that, in general, much more noise (other neurons are activated/deactivated) is introduced.

**Interpreting Metrics.**   While the reported values for precision and recall may appear low, we argue that these fragments are of good quality and point out that precision and recall metrics have to be interpreted differently than in a classification context: A net fragment consisting of many neurons characterizes (a part of) an object and maintains this characterization even when many of its neurons are inversed (and precision and recall are low). This aligns with the findings of Ahmad and Hawkins (2015), highlighting that binary distributed activations demonstrate high robustness, retaining accurate interpretation even in the presence of numerous flipped neurons.

In the case of the depicted straight lines, the feature neuron activations of the two most similar lines overlap with $21.2\%$. Therefore, achieving precision and recall values surpassing this threshold permits reliable discrimination even among similar lines, which is the case for the CNA but not for the autoencoder (e.g., see Figure 10).
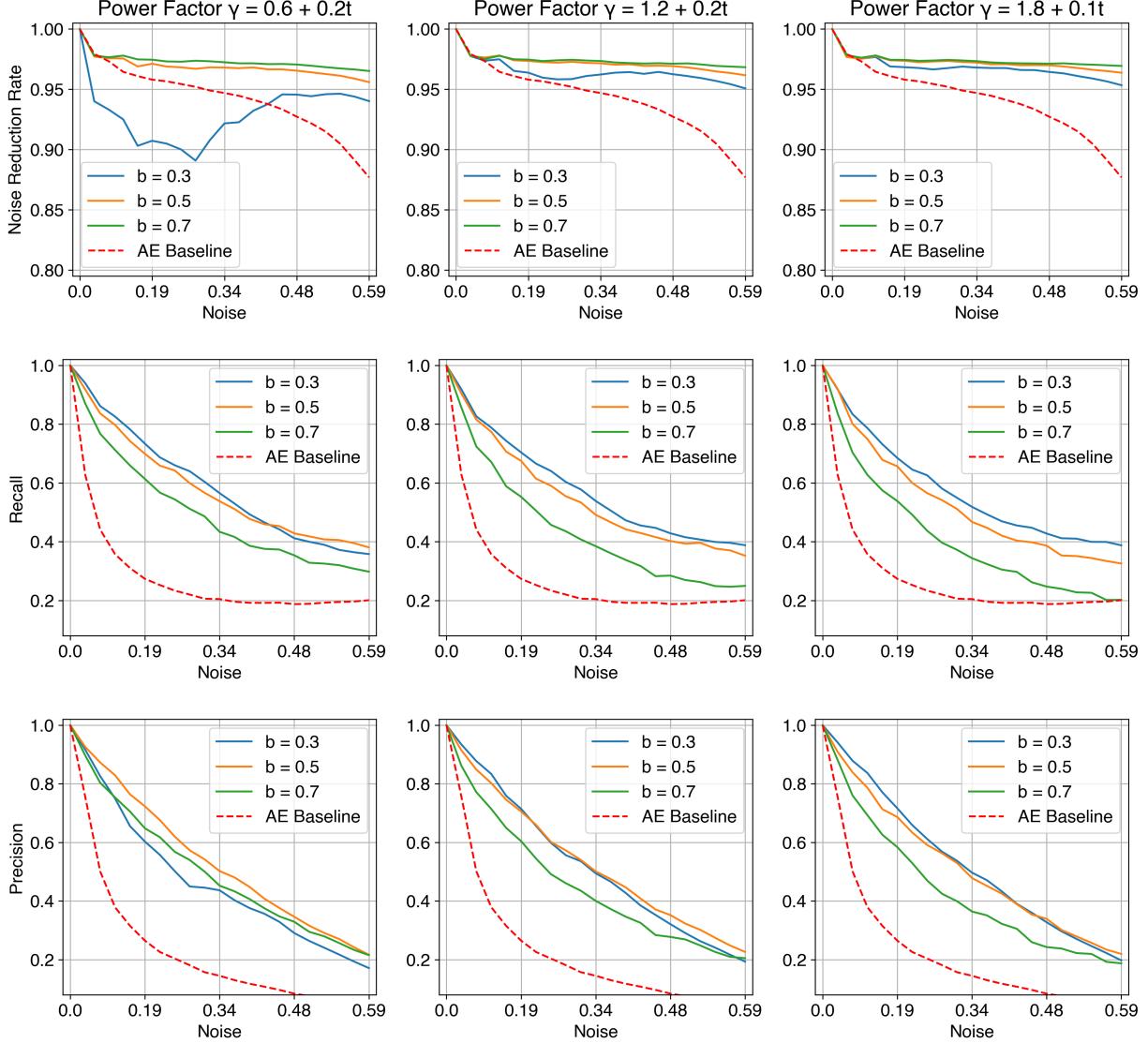
Figure 10: System performance as function of the noise level. First row: noise reduction rate; second row: recall; third row: precision. Columns differ in the attenuation coefficient $\gamma$. The colors within each plot represent different activation biases $b^{(S2)}$. The dotted red line is the autoencoder, serving as baseline.
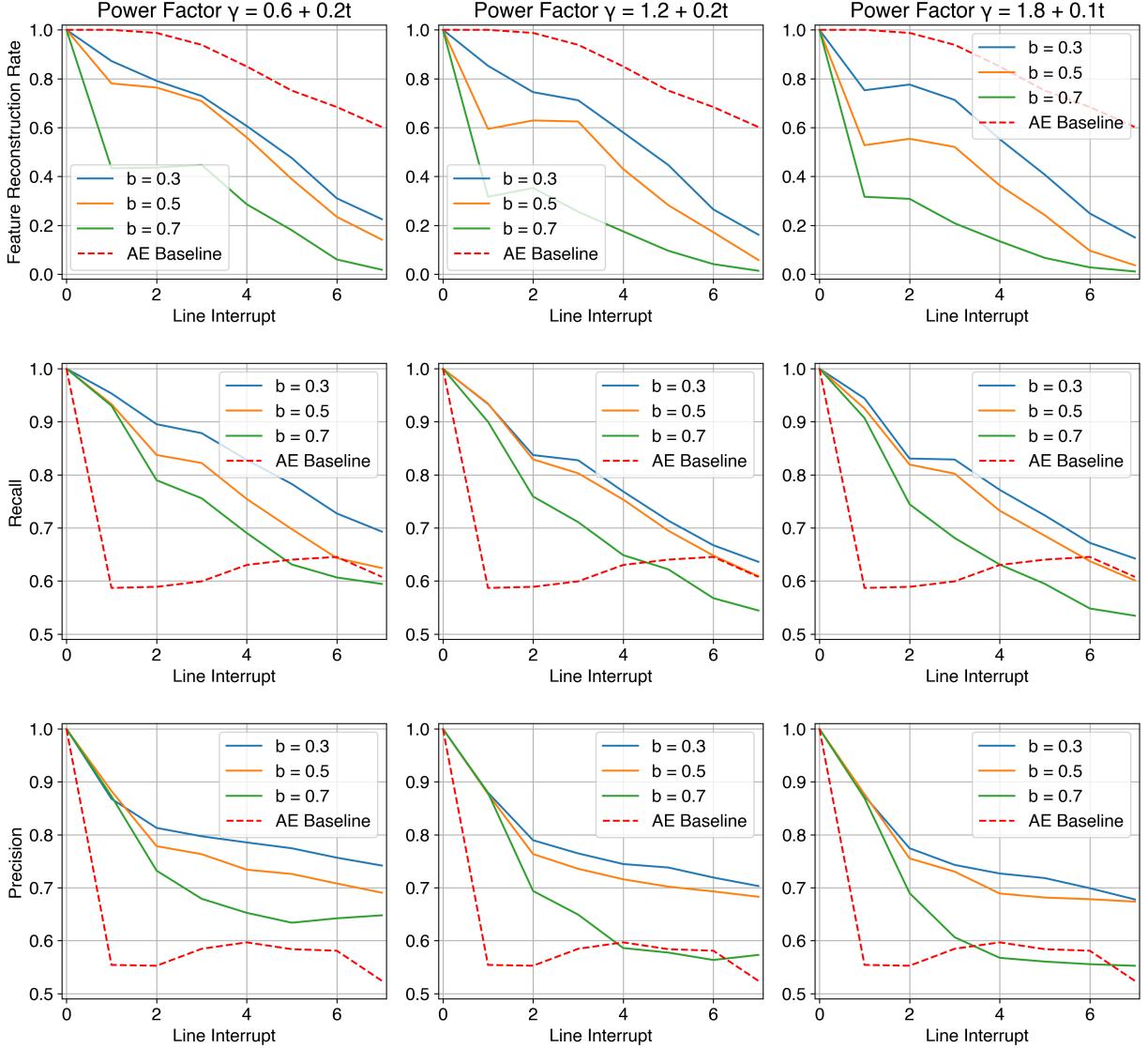
Figure 11: Reconstruction of partial occlusion. *First row:* feature reconstruction rate; *second row:* recall; *third row:* precision. Columns are distinguished by the attenuation coefficient $\gamma$, and line colors within plots distinguish activation bias $b^{(S2)}$. The dotted red line is the autoencoder model.

# References

Ahmad, S. and Hawkins, J. (2015). Properties of Sparse Distributed Representations and their Application to Hierarchical Temporal Memory. *arXiv*, (preprint arXiv:1503.07469).

Amari, S.-I. (1972). Learning Patterns and Pattern Sequences by Self-Organizing Nets of Threshold Elements. *IEEE Transactions on Computers*, C-21(11):1197–1206.

Amirian, M., Schwenker, F., and Stadelmann, T. (2018). *Trace and Detect Adversarial Attacks on CNNs Using Feature Response Maps*, volume 11081 of *Lecture Notes in Computer Science*, page 346–358.

Anderson, C. H. and van Essen, D. C. (1987). Shifter circuits: a computational strategy for dynamic aspects of visual processing. *Proceedings of the National Academy of Sciences*, 84(17):6297–6301.

Arathorn, D. W., Stevenson, S. B., Yang, Q., et al. (2013). How the unstable eye sees a stable and moving world. *Journal of Vision*, 13(10):22–22.

Carlucci, F. M., Russo, P., Tommasi, T., et al. (2019). Hallucinating Agnostic Images to Generalize Across Domains. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshop*, pages 3227–3234.

Cossell, L., Iacaruso, M. F., Muir, D. R., et al. (2015). Functional organization of excitatory synaptic strength in primary visual cortex. *Nature*, 518(7539):399–403.

Csurka, G. (2017). *Domain Adaptation in Computer Vision Applications*. Advances in Computer Vision and Pattern Recognition.

Hebb, D. O. (1949). *The Organization of Behavior; A Neuropsychological Theory*.

Hinton, G. and Salakhutdinov, R. (2006). Reducing the Dimensionality of Data with Neural Networks. *Science*, 313:504–7.

Hinton, G. E. and Sejnowski, T. J. (1983). Optimal Perceptual Inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.

Hopfield, J. J. (1982). Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences*, 79(8):2554–2558.

Kingma, D. and Ba, J. (2015). Adam: A Method for Stochastic Optimization. In *International Conference on Learning Representations (ICLR)*, San Diego, CA, USA.

LeCun, Y., Bengio, Y., and Hinton, G. (2015). Deep Learning. *Nature*, 521(7553):436–444.

LeCun, Y., Boser, B., Denker, J. S., et al. (1989). Backpropagation Applied to Handwritten Zip Code Recognition. *Neural Computation*, 1(4):541–551.

Memisevic, R. and Hinton, G. E. (2010). Learning to Represent Spatial Transformations with Factored Higher-Order Boltzmann Machines. *Neural Computation*, 22(6):1473–1492.

Moosavi-Dezfooli, S.-M., Fawzi, A., and Frossard, P. (2016). DeepFool: A Simple and Accurate Method to Fool Deep Neural Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2574–2582.

Mumford, D. (2002). Pattern Theory: the Mathematics of Perception. In *International Congress of Mathematicians*, volume 3, pages 1–21.

Olshausen, B., CH, A., and Van Essen, D. (1995). A Multiscale Dynamic Routing Circuit for Forming Size- and Position-Invariant Object Representations. *Journal of Computational Neuroscience*, 2:45–62.

Olshausen, B. A. and Field, D. J. (2005). How Close Are We to Understanding V1? *Neural Computation*, 17(8):1665–1699.

Palm, G. (2013). Neural associative memories and sparse coding. *Neural Networks*, 37:165–171.

Prince, S. J. D. (2023). *Understanding Deep Learning*.

Qi, F., Mattia, S., Yu-Win, g. T., et al. (2023). Towards Robust Object Detection Invariant to Real-World Domain Shifts. In *International Conference on Learning Representations (ICLR)*.

Rosenblatt, F. (1962). *Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms*.

Sager, P., Salzmann, S., Burn, F., et al. (2022). Unsupervised Domain Adaptation for Vertebrae Detection and Identification in 3D CT Volumes Using a Domain Sanity Loss. *Journal of Imaging*, 8(8):222.

Sager, P. J., Meyer, B., Yan, P., et al. (2025). AI Agents for Computer Use: A Review of Instruction-based Computer Control, GUI Automation, and Operator Assistants. *arXiv*, (preprint arXiv:2501.16150).

Simmler, N., Sager, P., Andermatt, P., et al. (2021). A Survey of Un-, Weakly-, and Semi-Supervised Learning Methods for Noisy, Missing and Partial Labels in Industrial Vision Applications. In *Proceedings of the 8th Swiss Conference on Data Science*, pages 26–31.

Tsodyks, M. and Feigel'man, M. (2007). The Enhanced Storage Capacity in Neural Networks with Low Activity Level. *EPL (Europhysics Letters)*, 6:101.

Tuggener, L., Emberger, R., Ghosh, A., et al. (2024). Real World Music Object Recognition. *Transactions of the International Society for Music Information Retrieval*, 7(1):1–14.

Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). Attention is All you Need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, volume 30, page 6000–6010.

von der Malsburg, C. (1981). The correlation theory of brain function. Internal Report 81-2, Max-Planck-Institut für Biophysikalische Chemie.

von der Malsburg, C. (2018). Concerning the Neuronal Code. *Journal of Cognitive Science*, 19(4):511–550.

von der Malsburg, C., Stadelmann, T., and Grewe, B. F. (2022). A Theory of Natural Intelligence. *arXiv*, (preprint arXiv:2205.00002).

Wilson, H. R. and Cowan, J. D. (1972). Excitatory and Inhibitory Interactions in Localized Populations of Model Neurons. *Biophysical Journal*, 12(1):1–24.

Wiskott, L., Fellous, J.-M., Krüger, N., et al. (1997). Face recognition by elastic bunch graph matching. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(7):775–779.

Wolfrum, P., Wolff, C., Lücke, J., et al. (2008). A recurrent dynamic model for correspondence-based face recognition. *Journal of Vision*, 8(7):34.