

# TOWARDS HUMAN-LIKE MACHINE VISION: REPRESENTING PART-WHOLE RELATIONSHIP WITH HIERARCHICALLY CORRELATED NEURONAL ACTIVATION IN NEURAL NETWORKS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Representing hierarchical structure is a key problem that characterizes the gap between current neural network and human-like intelligence. While human brain flexibly extracts part-whole hierarchy from unstructured sensory input, how can a neural network with fixed connection weight flexibly capture such compositional structure is still an open question. Most efforts in machine learning field focus on slot-based methods to temporally tackle the problem. In this paper, we provide new insights on this challenge without resort to the “slot” idea. From an interdisciplinary viewpoint that combine neuroscientific hypothesis and machine learning models, we propose the Composer, which dynamically “correlates” its distributed neural activation into an emergent implicit hierarchical structure to represent the part-whole hierarchy of objects. The observed representation is consistent to the widely-discussed “neural syntax” in neuroscience. Therefore, we hope the Composer shed light on a new paradigm to develop human-like vision and to build up compositional structure without “slots”. We also invent quantitative measures to evaluate the parsing quality, which shows that the Composer can parse a range of synthetic scenes of different complexities. By incorporating advanced machine learning models like LLMs or diffusion models into the paradigm, the capability of Composer is promising to be scaled into real-world datasets in the future. Taken together, we believe the Composer can inspire and inform future innovations and development towards artificial general intelligence (AGI).

## 1 INTRODUCTION

How symbol-like entities can emerge from distributed representations of neural networks and be composed into compositional structures are fundamental challenges underlining the gap between current neural network models and human-level artificial general intelligence (AGI) Greff et al. (2020). For example, infants as young as five months begin to make sense of objects Acredolo & Goodwyn (1988) but a general solution for object-centric representation in neural networks is still ongoing challenges. Greff et al. (2020) argues that the inability of existing neural networks to dynamically and flexibly “bind” information that is distributed throughout the network is the underlining cause of the discrepancy between their remarkable successes and chronic shortages, which still characterize a gap between deep learning and AGI.

Further, the Hinton (2021) continues to argue that the object-centric representation should also capture the part-whole relationship, which put higher demand of the compositionality of distributed representations in neural networks. In other words, to understand a visual scene, a concept, or a sentence, it should not only form single-level entities but a hierarchy of entities, and the “hierarchical relation” among multi-level entities should somehow also be represented within a neural network. Since both the identification of object itself and the relation among objects are likely to be uncertain and diverse, how can a pure neural network with fixed architecture and connection capture such flexibility? On the other hand, the ability to form such hierarchical representation is highly desirable because it could potentially facilitate planning efficiency and hierarchical planning on complex problems can be translated into the part-whole representation problem Rao et al. (2022).

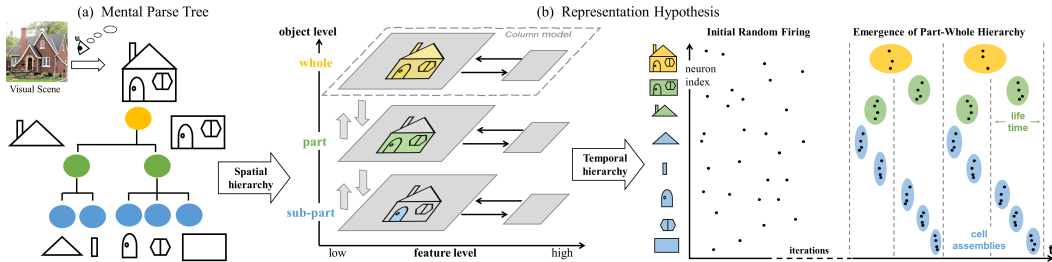


Figure 1: (a) The mental parse tree of the visual scene. (b) The neural representation of the mental parse tree. Left: Spatial separation of neural space into levels. Colored areas indicate activated neurons. Right: Represent the parse tree as nested neuronal coherence. The y-axis of spike raster plot are neuron index based on spatial hierarchy. Correlated cell assemblies are indicated by colored shadows.

Most efforts to address this challenge in machine learning center on slot-based methods. These approaches firstly divide their latent representations into pre-defined “slots”, creating a discrete separation of object representations. These slots are easier to be further organized into compositional structures, eg. in graph neural network Han et al. (2022); Xu et al. (2017), generative models Deng et al. (2021), attention-based models Sun et al. (2021); Fisher & Rao (2022) or capsule networks Hinton et al. (2018); Garau et al. (2022). However, this simplicity may limit slot-based approaches from representing the full diversity of objects and relationship Lowe et al. (2023). Their pre-defined and discrete nature in principle make it challenging to account for the flexibility and uncertainty of the representation; and it seems unlikely that the brain assigns entirely distinct groups of neurons to each perceived object.

The motivation of this paper is to develop a framework to represent part-whole hierarchy without discrete slots, but with continuous and distributed object-centric representations. More specifically, distributed neural activation is dynamically “correlated” into different “groups” of neurons, each of which have more coherent inner-group patterns. This idea is inspired by the neural syntax hypothesis Buzsáki (2010) in neuroscience: The distributed neuronal firings are dynamically correlated into cell assemblies (groups of correlated neurons), which is believed to be the neural words to construct neural syntax in the brain. Further, the nestedness of the hierarchical relationship is represented as dynamically formed nestedness of neuronal coherence, or spatial-temporally nested cell assemblies (Fig.1b,right).

While the dynamical and continuous nature of neuronal coherence fits the demands of diversity, flexibility, and uncertainty, how the nested states can emerge seems a “magic”. To unfold the magic, we take inspirations from a recently developed cortical-column-inspired machine learning model Zheng et al. (2022), where correlated cell assemblies dynamically emerge to form object-centric representations. Taking this model as the building block, we developed the **Composer** (short for **C**Ortical-inspired **e**Mergence of **P**art-wh**O**le relation**S**hip through **n**Eurolal cohe**R**ence), to further study how the ‘neural words’ are composed into hierarchical ‘neural syntax’.

The last challenge is how to explicitly evaluate the nestedness of the part-whole representation since the correlation patterns are diverse, uncertain, distributed and of hierarchical levels. To this aim, we developed four synthetic datasets of different complexities and three quantitative metrics to measure different aspects of the part-whole representation. Quantitative results and qualitative visualization confirm the validity of the Composer and the plausibility of the framework.

## 2 NEUROSCIENTIFIC MOTIVATION

A widely discussed hypothesis in neuroscience is that transiently active ensembles of neurons, known as the “cell assembly,” represents a distinct cognitive entity Hebb (2005). The cell assemblies are dynamically correlated and organized on gamma rhythms, which reflects a coherent state. It posits that the oscillating firing patterns of neurons give rise to two message types: the discharge frequency encodes the presence of features, and the relative timing of spikes encode feature binding Singer

(2009); Malsburg (1994). As a result, the temporally correlated spike firings ‘bind’ distributed features into an entity or neural words. Further, it is hypothesized that the hierarchical organization of assembly sequence may be regarded as a neural syntax, so that the hierarchical relationship is represented as the spatial-temporally nested structure of cell assembly trajectories, accompanied by nested neural rhythms Buzsáki (2010). Such nested states is argued to be the pre-requisite of consciousness Northoff & Zilio (2022). Taken together, these hypothesis motivated a coherence-based representation paradigm for the part-whole hierarchy.

It is further discussed that the formation and readout of coherent cell assemblies depend on several mechanisms. Firstly, the time-window of the readout neurons should fit the time-scale of cell assemblies and there should be a hierarchical organization of time-windows Buzsáki (2010). Secondly, the neural oscillation is closely related to the delay-coupled top-down feedback, which provides predictive modulations Singer (2021). Such delayed-coupled structure is considered as a basic feature of cerebral cortex. Thirdly, the oscillatory dynamics is partially determined by the pre-configured wiring in the brain Buzsáki & Mizuseki (2014). As unfolded in Section 5, these mechanism motivated the architecture of the Composer.

### 3 REPRESENTATION HYPOTHESIS

The first challenge to represent part-whole hierarchy is the co-existence of multi-level objects in the same representation space. Since part-whole level is distinguished in nature from conventional layers in a deep network, we distinguish two types of ‘levels’ to avoid confusion (Fig.1b,left). One is feature level, corresponding to layers in deep networks. Taking autoencoder as example, neurons at first layer encode pixel-level features while neurons at latent space encode more abstracted features, yet of the same object. Another is object level, which forms neuronal coherence at different spatial-temporal scales to represent objects of distinct part-whole levels. All object levels share the same set of feature levels. For example, the low-level features of each object level is pixel-level features, but correlated into different coherent firing states so as to represent parts or wholes (colored areas in Fig.1b, left). One analogy is that feature levels are different layers in a cortical column while object levels are different cortical regions, sharing similar 6-layer column structures (Fig.2g). At each object level, we assume the firing of spiking neurons are dynamically synchronized into cell assemblies to bind the features encoding the object at that level. Multiple part objects are represented as a sequence of cell assemblies. Assemblies at higher-object-levels have larger spatial size and longer life time than that at lower-object-levels (Fig.1b).

The second challenge is to represent the hierarchical relationship. We hypothesize that the longer life time of each whole-level cell assembly covers the sequence of its part-level assemblies. As a result, the nested structure in a part-whole relationship is represented as the nested temporal structure of cell assembly sequences (Fig.1b, right). If we consider pixel-level assemblies, they are also spatially nested since pixels of whole-level object covers that of its parts (colored areas in Fig.1b, left). The analogy of shared spatial map (pixels) across object levels is the topographical mapping in the cortical organization Kaas (1997).

### 4 THE MAGIC

Given the representation hypothesis, it is still a mystery how such specific coherent state can emerge given a visual scene? Are there unknown magic to carefully organize neurons into assemblies and align them nicely into nested states? Recently, Zheng et al. (2022) proposed DASBE (Fig.2c,d), a cortical-column inspired hybrid model to generate oscillatory single-level cell assemblies as object-centric representation through the iterative bottom-up / top-down processing between a spike coding space (SCS) and a pre-trained denoising autoencoder (DAE). The insight is two-fold. Firstly, it was proved that the denoising autoencoder can parameterize a recurrent dynamics to build up many associative memories or attractive states:  $x_{t+1} = DAE(x_t)$ . For example, each object in the training set becomes a memory of the dynamics (Fig.2a,b). Given an multi-object image, the retrieval of the memory actually activates the neurons encoding one of the objects. Secondly, the refractory period of spiking neurons “destroy” the attracted states so that the activation becomes transient. Thus, alternative groups of neurons compete with each other to get activated (by DAE) and then get silenced

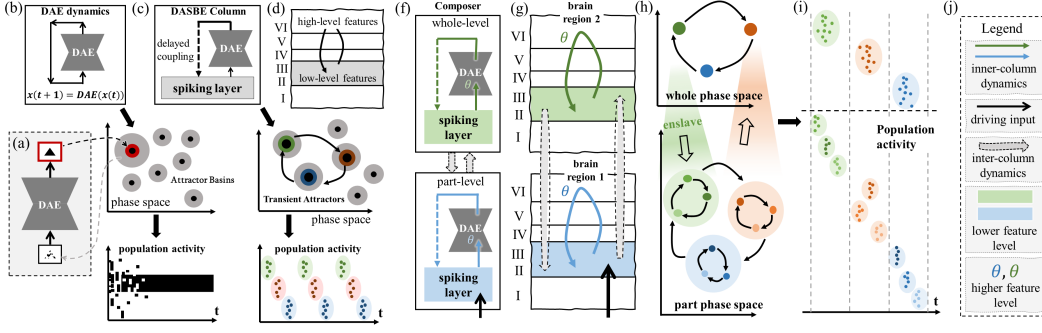


Figure 2: How nested neuronal coherence emerges. (a)~(d): Insights from DASBE Zheng et al. (2022). (f)~(j): Insights of Composer. (a) Denoising autoencoder (DAE). (b) Build up attractor dynamics by DAE (top and middle), which results in stationary population activity (bottom). (c) DASBE column: Break attractors to induce rhythmic dynamics by spiking neurons and delay coupling (middle and bottom). (d) Analogy between DASBE and single cortical column. (e) General architecture of the Composer and (g) the analogy with multi-level columns (h) The “enslave” dynamics of Composer and (i) resulting nested population activity. (j) Legend for (f)(g).

(by SCS). The DAE feedback delay and integration time window softly shape the temporal structure of the assembly sequence.

Inspired by the insight, we treat each DASBE as a ‘cortical column’ of a single object level, where the pixel-wise SCS and the latent space of DAE correspond to low and high feature levels (Fig.2f,g). To represent objects of hierarchical levels, several DASBE columns are “stacked” and interconnected. The DASBE column at whole levels have longer integration time window than that of part-levels so that the assembly have longer life time. Besides, the prior of part / whole objects are introduced to DAEs by pre-training. Lastly, whole-level assemblies “enslave” the dynamics of part-level assembly trajectory by gating mechanism (Fig.2h,i). Taken together, delayed-coupled column architecture, time window hierarchy, pre-configured wiring and predictive gating are core mechanisms to unfold the magic.

## 5 MODEL

In this work, we consider a whole level and a part level and only focus on the neuronal coherence of pixel-level features to simplify the problem. Such setting can be generalized in future works (See A.1). Overall, the Composer consists of two levels of columns, interconnected by top-down modulation and bottom-up integration (Fig.2f). Each column contains a pixel-level spiking layer, named spike coding space (SCS, Fig.3c) and a denoising autoencoder (DAE). ‘Pixel-level’ implies that SCS of both levels have the same dimensions as the image (d). SCS and DAE are delay-coupled in each column. More specifically, spikes in SCS are integrated within a narrow time window, known as coincidence detectors König et al. (1996) and the delayed feedback from DAE modulates the firing rate of spiking neurons by gating, consistent with dendritic computation of pyramidal cells in cortex Sherman & Guillery (1998).

### 5.1 PART-LEVEL COLUMN

Spiking neurons in the  $SCS_1$  of part column receive inputs from three sources (Fig.3b): the input image  $x \in \{0, 1\}^d$ , the inner-column feedback  $\gamma_1(t) \in \mathbb{R}^d$  and the inter-column feedback  $\Gamma(t) \in \mathbb{R}^d, t \in [0, T]$ :

$$\rho_1(t) = x \cdot \gamma_1(t) \cdot \Gamma(t) \quad (1)$$

where ‘ $\cdot$ ’ is pixel-wise and  $\rho_1$  is the firing rate, which determines the firing activity  $s_1 \in \{0, 1\}^d$ :

$$P(s_1 = 1) = \rho_1(t) \cdot g_1(t - \hat{t}), \quad t \in [0, T]. \quad (2)$$

where  $g_1(t - \hat{t})$  is the relative refractory function of neurons in the part level (Fig.3e) and  $\hat{t}$  is the timing of the latest spike firing event of each neuron. As shown in Fig.3e, after firing a spike, the neuron goes into an absolute refractory period of timescale  $\tau_{r1}$  and then a relative refractory period of timescale  $\tau_{\delta 1} - \tau_{r1}$ , where firing probability is inhibited by a factor  $g < 1$ . The inner-column feedback  $\gamma_1$  in eq.1 is the denoised output of the  $DAE_1$ , yet with delay  $\tau_d$ :

$$\gamma_1 = DAE_1([I_1 * s_1](t - \tau_d)). \quad (3)$$

where  $*$  is the convolution operator and  $I_1$  is the integrative function for  $s_1(t)$ , of time window  $\tau_1$  (Fig.3f):  $[I_1 * s_1](t) = \sum_{\tau=0}^{\tau_1} (I_1(\tau) \cdot s_1(t - \tau))$ . In a word, the spiking activity  $s_1(t)$  in the SCS is integrated within a short time window  $\tau_1$  before fed into the  $DAE_1$ , and the feedback from  $DAE_1$  to  $SCS_1$  is delayed by  $\tau_d$ .

## 5.2 WHOLE-LEVEL COLUMN

Since the whole-level column is the top level in the current two-level Composer, it does not receive top-down modulation from even higher levels. The firing rate of spiking neurons in  $SCS_2$  is determined by the multiplication of two sources: a driving term and a modulation term.

$$\rho_2 = (\lambda \cdot x + (1 - \lambda) \cdot D) \cdot \gamma_2 \quad (4)$$

where  $\lambda < 1$  is the factor of partial influence from skip connection (Fig.3c).  $D$  is the integrated spikes from the part-level  $SCS_1$ , serving as the main driving input for the whole-level column.  $\rho_2$  similarly determines the spike firing in  $SCS_2$  by:

$$P(s_2 = 1) = \rho_2(t) \cdot g_2(t - \hat{t}) \quad (5)$$

Besides,  $\gamma_2$  in eq.4 is the delayed feedback from  $DAE_2$  (similar to eq.3):

$$\gamma_2 = DAE_2([I_2 * s_2](t - \tau_d)) \quad (6)$$

## 5.3 LINKING THE LEVELS

Up to now, we have introduced operations within each column of the Composer except for two variables:  $\Gamma$  in eq.1 and  $D$  in eq.4, which are interactions between levels:

$$\Gamma(t) = [I_\Gamma * s_2](t - \tau_{d'}) \quad \text{and} \quad D(t) = [I_D * s_1](t). \quad (7)$$

where  $I_\Gamma, I_D$  is the corresponding integration function.  $\tau_{d'}$  is the cross-level feedback delay, which is set to be  $\tau_{d'} = \tau_d$ .

It is notable that each integration function  $I_i$  has its time window parameter  $\tau_i$ , ( $i \in \{1, 2, \Gamma, D\}$ ) and each refractory function  $g_i$  has absolute refractory period  $\tau_{ri}$  within the entire refractory period  $\tau_{\delta i}$  ( $i \in \{1, 2\}$ ). The general relationship is that  $\tau_1 < \tau_2 < \tau_D \sim \tau_\Gamma \ll \tau_d$  and  $\tau_{\delta 1} \sim \tau_{\delta 2} \ll \tau_d$ . As motivated by neuroscience, the hierarchical organization of integration time windows and delay coupling are essential mechanisms for cell assemblies. These time-scale parameters are treated as hyper-parameters. See A.4.2 for more details and for ablation studies. In general, the performance of the model is robust to variation of timescale parameters within a reasonable range. It is also notable that the  $DAE_1$  and  $DAE_2$  are pre-trained to denoise part objects and whole objects respectively, as the pre-configured wiring in the cortical columns. See A.5 for details of training.

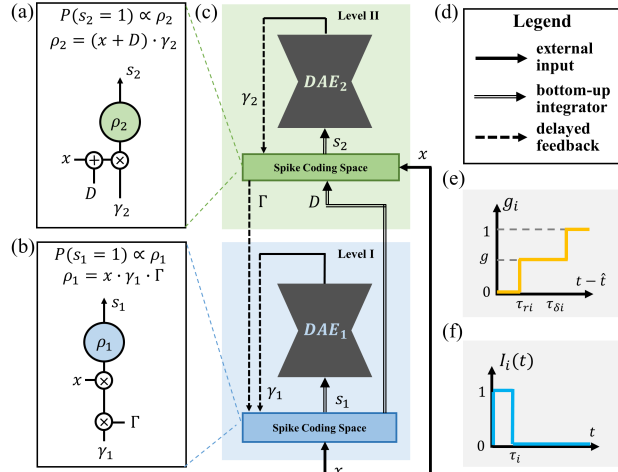


Figure 3: (a)(b) Pyramidal neuron models in the spike coding space of part level and whole level.  $\otimes$  stands for multiplication and  $\oplus$  stands for addition on the dendrites. (c) Detailed information flow. Note that delayed coupling exists both within each column and between different columns. Levels are indicated by color. (d) The legend for (c). (e) Relative refractory function  $g_i$ . (f) Integration function  $I_i(t)$  of timescale  $\tau_i$ .

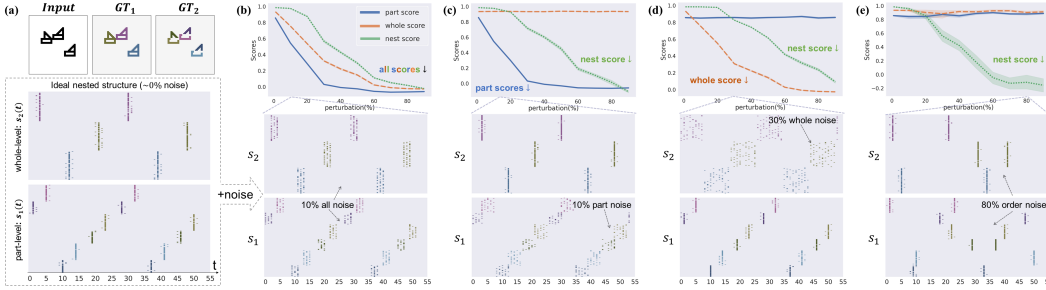


Figure 5: Intuitions of metrics. (a) Ideal nested structure (spike raster plot of two levels) given the input (top). If random noises are added to (c) part-level, (d) whole-level, or (b) both levels, corresponding scores degrades gradually from 1 to 0. If correlated noise are added to perturb the hierarchical relation among part assemblies and whole assemblies, only nest score degrades. (b)~(e): Top, scores vs perturbation degree in each case; Bottom, exemplified perturbed spike pattern of certain degree. See Fig.11~Fig.14 for details.

## 6 EVALUATION

Quantitative evaluation of the hierarchical structure is essential to verify or falsify models, yet is missing in most related works Hinton et al. (2018); Garau et al. (2022); Sun et al. (2021). The challenge comes from two parts: (1) real-world images do not have explicit parsing structures and (2) it is non-trivial to measure distributed hierarchical structures. These challenges motivate us to invent synthetic datasets and metrics to explicitly evaluate the Composer.

### 6.1 DATASET

We invent four synthetic part-whole datasets of different complexities (Fig.4), each containing 60000 samples. Each image consists of multiple whole objects. Each whole object is further composed of well-defined parts.

**Ts** dataset (Fig.4a) consists of 3 letter  $\top$  and 3 reversed letter  $\perp$  as whole-level objects. Each  $\top$  or  $\perp$  is composed of a horizontal bar and a vertical bar as parts. *T*'s dataset has larger whole number. Similar stimuli are used as target templates in perceptual tasks in neuroscience Wolfe (2021).

**Squares** dataset (Fig.4c) contains 3 randomly-located squares as wholes, each of which consists of 4 corners. This dataset has more parts. Gestalt psychology applies similar stimuli to study illusory contour Lee & Nguyen (2001).

**SHOPs** (short for **S**hoes (ii), **H**ouse (i), **O**pera (iii) in Fig.4b) consists of 3 types of wholes that are further composed of more elementary rectangular and triangles. This dataset accounts for the complexity that parts could overlap to construct a whole. Overlapped pixels are not assigned to either part in the ground truth (Fig.4b, bottom).

**Double-Digit MNIST** (Fig.4d) contains 2 randomly selected and located double digits, and each double-digit is composed of 2 closely located MNIST digits. This dataset contains objects of higher complexity and diversity.

### 6.2 SCORES

How to evaluate the nestedness in Fig.1b (right)? The idea is that the correlated cell assemblies can be translated into clusters of spike trains. Therefore, the coherence measure for clustering

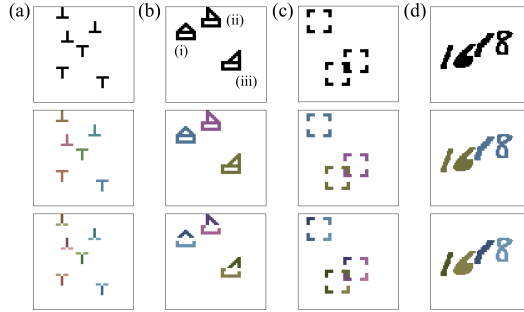


Figure 4: Exemplified samples in datasets: (a) Ts (b) SHOPs (c) Squares (d) Double-Digit MNIST. Top, input. Middle / Bottom, ground truth of Wholes / Parts.

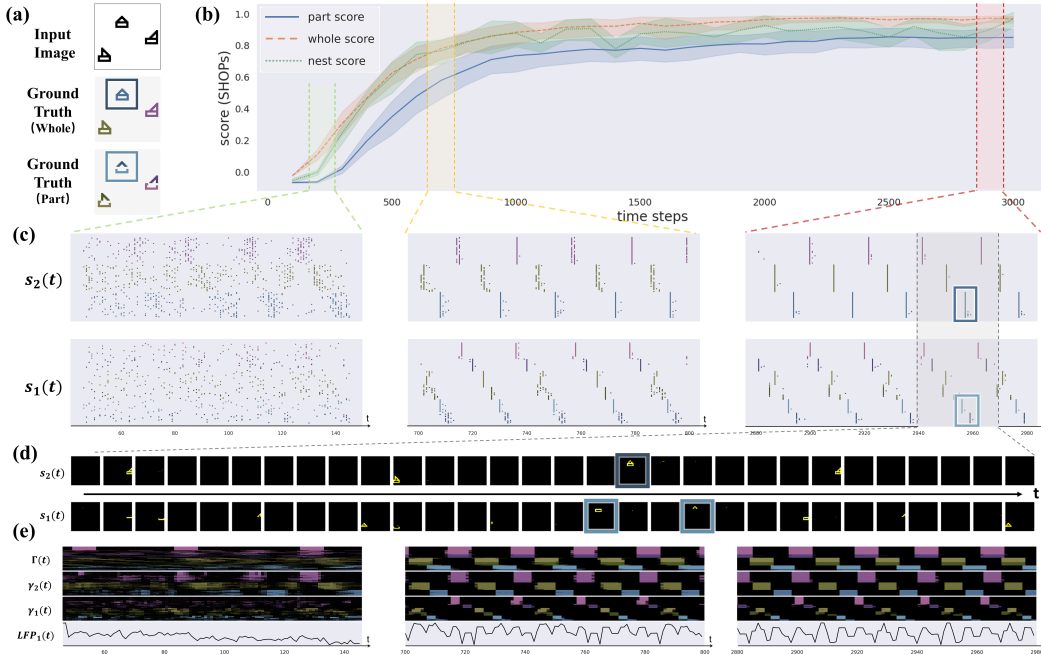


Figure 6: Emergence of the part-whole hierarchy with nested neuronal coherence. Exemplified by one SHOPs sample. (a) Input image and ground truth; (b) Evolution of the Scores. (c) The spike raster plot of three selected phases in (b): phase I (initial, green box), phase II (middle, yellow box), phase III (final, red box).  $s_2(t)$ ,  $s_1(t)$  stand for spiking representations in SCSs of whole/part levels. (d) zoomed in spiking pattern during the period marked by black box in (c), to visualize what each cell assembly represents (e.g. blue boxes in (c),(d)and(a)). (e) Evolution of the top-down attention maps and the local field potential (LFP) during the three phases in (b).

method, like Silhouette Score Rousseeuw (1987), becomes the candidate. On this basis, we develop Part Score =  $\text{Silhouette}(s_1, s_1, GT_1, D_{vp})$  and Whole Score =  $\text{Silhouette}(s_2, s_2, GT_2, D_{vp})$  to measure the neuronal coherence within each level.  $s_{1/2}$  is the collection of spike trains and  $GT_{1/2}$  is the ground truth of spike trains at each level (Fig.5a).  $D_{vp}$  is the non-Euclidean Victor-Purpura metric Victor & Purpura (1996) for evaluating distance among spike trains (See A.3.1).  $\text{Silhouette}(a, b, c, d)$  measures the inner-cluster coherence among sample set  $a$  and sample set  $b$  (generally  $a = b$ ), given cluster assignment  $c$  and distance metric  $d$ . To evaluate the nested structure, we need to generalize the Silhouette to account for cross-level coherence. Therefore, we define Nest Score  $\propto \text{Silhouette}(s_1, s_2, GT_2, D_{vp})$ , which measures coherence between part-level spike trains  $s_1$  and whole-level spike trains  $s_2$  given whole-level cluster assignment  $GT_2$  (also see A.3). It is notable that the shared pixel-wise spatial map of  $s_1, s_2, GT_1, GT_2$  guarantees their comparability.

To show the intuition of the scores, we gradually perturb an ideal nested structure (Fig.5a) into noisy patterns and see how scores change. If only the part level is perturbed, the Whole Score remains constant while both the Part Score and Nest Score are decreased from 1 to 0 (Fig.5c), vice versa (Fig.5d). If we only perturb the relative order of part assemblies and whole assemblies to degrade the nestedness, then the Nest Score decreases saliently while Part and Whole Scores remain (Fig.5e). Taken together, three scores completely capture the structure of nested neuronal coherence. See A.3.5 for more details. Following the Silhouette Score, the best score (coherence) is 1 and the worst score (incoherence) is -1. A score near 0 indicates randomness.

## 7 EXPERIMENTS

### 7.1 QUALITATIVE RESULTS AND VISUALIZATION

**Emergence of the “neural syntax” in SCS.** We visualize the simulation on a randomly selected sample in the SHOPs dataset in Fig.6. As indicated by the convergence of Part Score, Whole

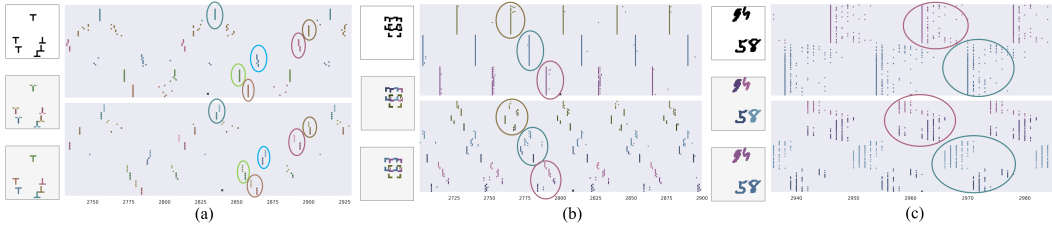


Figure 7: Visualization of the emergent part-whole hierarchy on other datasets: (a) Ts (b) Squares (c) Double-Digit MNIST. (In each sub-figure) Left: input image, part-level ground truth and whole-level ground truth. Right: spike raster plot in final phase III, top for whole level and bottom for part level. Cell assemblies are circled for clarification. See Fig.19~Fig.26 for complete visualizations

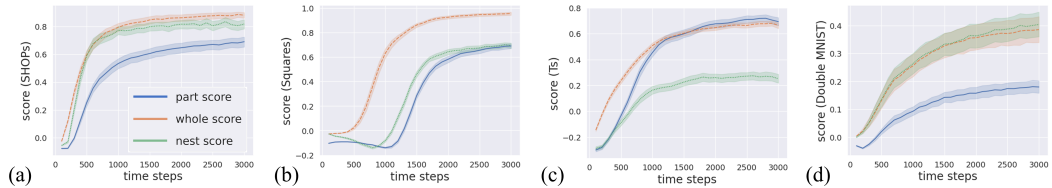


Figure 8: Convergence of Scores. (a) SHOPs (b) Squares (c) Ts (d) Double-Digit MNIST.

Score and Nest Score (Fig.6b), the Composer gradually achieves a state of neuronal coherence that represents the parts and wholes as synchronized cell assemblies, which is further visualized in Fig.6c right and Fig.6d. More specifically, three two-level binary-tree (corresponding to three SHOPS objects in Fig.6a) periodically emerges in the final phase III (Fig.6c right), one of which is marked out by deep blue (for whole object) and light blue (for part objects) boxes. The spikes for part / whole objects in Fig.6c are reordered (on the y-axis) and colored corresponding to the ground truth of part /whole objects (Fig.6a) for more vivid visualization, so that the same cell assemblies are arranged closely and the color of their spikes is consistent with the ground truth. For more visualization results, see A.7.3.

**Emergence of the hierarchical attention sequence** is observed in Fig.6e: Starting from randomness (left), the cross-level feedback  $\Gamma$  or inner-level feedback from the DAEs  $\gamma_1/\gamma_2$  gradually converge to a sequential trajectory. The similarity between assembly sequence (Fig.6c) and attention sequence (Fig.6e) indicates the close relationship between predictive top-down feedback and oscillatory cell assemblies Engel et al. (2001), consistent with the neuroscientific motivation. Therefore, the top-down attention from DAEs and bottom-up integration of spikes work together in Composer, leading to co-emergence of both attention sequence and neuronal coherence. Besides, rhythmic population activity ( $LF P_1$ ) emerges (Fig.6e).

**Visualization results on other datasets** are shown in Fig.7. Interestingly, the emergent coherence structure differs across the datasets. While in the Ts dataset (consists of 6 Ts), 6 binary trees emerge periodically, 3 quadrees emerge in the Squares dataset (consists of 3 Squares). In Double-Digit MNIST (consists of 2 Double-MNISTs), 2 binary trees emerge. Taken together, the Composer successfully and flexibly represents the part-whole hierarchy of scenes of different complexities.

## 7.2 QUANTITATIVE ANALYSIS

**Convergence** of the scores during iterations are evaluated on 100 randomly selected samples in each dataset and are shown in Fig.8. Interestingly, while scores consistently converge on all datasets with low error bars, the convergent process slightly differs across cases. For Squares (Fig.8b), whole objects group much faster than part objects, similar to human vision Lee & Nguyen (2001). For Ts (Fig.8c), the large object number imposes combinatorial burdens on the cross-level coordination, so that the Nest Score lags behind. For Double-Digit MNIST (Fig.8d), Composer has more difficulties in distinguishing part-level MNISTs, partially due to the diversity of the dataset. See A.5.5 for details.



## 8 RELATED WORK

Slot-based models are a line of research to represent object by pre-defined discrete slots Greff et al. (2015; 2016; 2017; 2019; 2020); Locatello et al. (2020). This idea is also applied to generate “scene graph” in graph neural networks Han et al. (2022); Xu et al. (2017), generative models Deng et al. (2021), attention-based models Sun et al. (2021); Fisher & Rao (2022) or capsule networks Hinton et al. (2018).

In this paper, we look for mechanisms without discrete and pre-defined slots, but with distributed and dynamical “coherence”. This idea has fewer related works and coherence was defined in diverse ways: the phase of complex-values Löwe et al. (2022), timing of spike firings Zheng et al. (2022), or rotation of vectors Lowe et al. (2023). But these works only consider single-level objects instead of part-whole relationship. The identical islands of vectors in GLOM and its variants Hinton (2021); Garau et al. (2022) are promising to account for part-whole via coherence of vectors. While GLOM exploit spatial coherence measure of vectors the Composer exploit both spatial and temporal coherence of neuronal activities.

## 9 DISCUSSION AND CONCLUSION

While the Composer is still a toy model that is being evaluated on toy datasets, we believe it can inspire and inform broader innovations to deal with more complex cases. It is notable that the design principle of the Composer is extremely general: it combines time scale hierarchy as biological constraint and denoising autoencoder as machine learning toolkit. So that the nested coherence states emerge through the iterative bottom-up / top-down processing between neuroscience-constraint SCS and machine-learning-based DAE. Since denoising is a fundamental operation in deep learning field, which is shared in a list of more advanced models like BERT Devlin et al. (2019), Diffusion model Yang et al. (2023), and LLMs and the denoising ability of the DAE is positively related to the parsing ability of the Composer (Fig.18), it is promising to incorporate the idea with alternative advanced models on diverse tasks. On the one hand, such representation framework contributes to interpretable representations; On the other hand, the idea of representing entity and relation through correlation or coherence instead of discrete slots is more likely to capture the feature of AGI and brain.

Also, in the Composer, the attention and coherence co-emerge along the simulation dynamics, and can be unified as associative memories (Fig.2c,h). Therefore, attention mechanism becomes another linking point with LLMs. Indeed, attention can be unified as associative memories, for example the self-attention in Transformer Vaswani et al. (2017) can be taken as the update function of a dynamical system Ramsauer et al. (2020). Therefore, it is intriguing to see how spatial-temporal “neural syntax” is realized in brain-inspired LLMs.

In conclusion, we look for solutions to represent compositional structure in neural networks without slots. To this aim, we combine neuroscience and machine learning to develop the Composer, which dynamically build-up “neural syntax” of “neural words” from distributed representation. We hope the Composer advances a new paradigm to develop human-like vision in the future. See A.1 for limitations and broader future works.

## REFERENCES

- Moshe Abeles. Role of the cortical neuron: integrator or coincidence detector? *Israel journal of medical sciences*, 18 1:83–92, 1982.
- Linda P. Acredolo and Susan W. Goodwyn. Symbolic gesturing in normal infants. *Child development*, 59 2:450–66, 1988. URL <https://api.semanticscholar.org/CorpusID:41352969>.
- André Moraes Bastos, W. Martin Usrey, Rick A Adams, George R. Mangun, Pascal Fries, and Karl J. Friston. Canonical microcircuits for predictive coding. *Neuron*, 76:695–711, 2012.
- Lars Buesing, Johannes Bill, Bernhard Nessler, and Wolfgang Maass. Neural dynamics as sampling: a model for stochastic computation in recurrent networks of spiking neurons. *PLoS computational biology*, 7(11):e1002211, 2011.

- György Buzsáki. Neural syntax: Cell assemblies, synapsembles, and readers. *Neuron*, 68:362–385, 2010.
- György Buzsáki. The brain from inside out. 2019.
- György Buzsáki and Kenji Mizuseki. The log-dynamic brain: how skewed distributions affect network operations. *Nature Reviews Neuroscience*, 15:264–278, 2014. URL <https://api.semanticscholar.org/CorpusID:18601814>.
- György Buzsáki, Adrien Peyrache, and John L. Kubie. Emergence of cognition from action. *Cold Spring Harbor symposia on quantitative biology*, 79:41–50, 2014. URL <https://api.semanticscholar.org/CorpusID:6613459>.
- Peter Dayan and L. F. Abbott. Theoretical neuroscience: Computational and mathematical modeling of neural systems. 2001.
- Fei Deng, Zhu Zhi, Donghun Lee, and Sungjin Ahn. Generative scene graph networks. In *International Conference on Learning Representations*, 2021.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics*, 2019. URL <https://api.semanticscholar.org/CorpusID:52967399>.
- Rodney J. Douglas and Kevan A. C. Martin. Neuronal circuits of the neocortex. *Annual review of neuroscience*, 27:419–51, 2004.
- Maria K. Eckstein, Christopher Summerfield, Nathaniel D. Daw, and Kevin J. Miller. Predictive and interpretable: Combining artificial neural networks and classic cognitive models to understand human learning and decision making. *bioRxiv*, 2023. URL <https://api.semanticscholar.org/CorpusID:258788055>.
- Simon B. Eickhoff, R. Todd Constable, and B. T. Thomas Yeo. Topographic organization of the cerebral cortex and brain cartography. *NeuroImage*, 170:332–347, 2017.
- Andreas Karl Engel and Wolf Singer. Temporal binding and the neural correlates of sensory awareness. *Trends in Cognitive Sciences*, 5:16–25, 2001.
- Andreas Karl Engel, Pascal Fries, and Wolf Singer. Dynamic predictions: Oscillations and synchrony in top-down processing. *Nature Reviews Neuroscience*, 2:704–716, 2001.
- Ares Fisher and Rajesh P. N. Rao. Recursive neural programs: Variational learning of image grammars and part-whole hierarchies. *ArXiv*, abs/2206.08462, 2022.
- Karl J. Friston. The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*, 11:127–138, 2010.
- Nicola Garau, Niccoló Bisagno, Zeno Sambauro, and Nicola Conci. Interpretable part-whole hierarchies and conceptual-semantic relationships in neural networks. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 13679–13688, 2022.
- Wulfram Gerstner, Werner M. Kistler, Richard Naud, and Liam Paninski. Neuronal dynamics: From single neurons to networks and models of cognition. 2014.
- Klaus Greff, Rupesh Kumar Srivastava, and Jürgen Schmidhuber. Binding via reconstruction clustering. *arXiv preprint arXiv:1511.06418*, 2015.
- Klaus Greff, Antti Rasmus, Mathias Berglund, Tele Hao, Harri Valpola, and Jürgen Schmidhuber. Tagger: Deep unsupervised perceptual grouping. *Advances in Neural Information Processing Systems*, 29, 2016.
- Klaus Greff, Sjoerd Van Steenkiste, and Jürgen Schmidhuber. Neural expectation maximization. *Advances in Neural Information Processing Systems*, 30, 2017.

- Klaus Greff, Raphaël Lopez Kaufman, Rishabh Kabra, Nick Watters, Christopher Burgess, Daniel Zoran, Loic Matthey, Matthew Botvinick, and Alexander Lerchner. Multi-object representation learning with iterative variational inference. In *International Conference on Machine Learning*, pp. 2424–2433. PMLR, 2019.
- Klaus Greff, Sjoerd Van Steenkiste, and Jürgen Schmidhuber. On the binding problem in artificial neural networks. *arXiv preprint arXiv:2012.05208*, 2020.
- Kai Han, Yunhe Wang, Jianyuan Guo, Yehui Tang, and Enhua Wu. Vision gnn: An image is worth graph of nodes. *ArXiv*, abs/2206.00272, 2022.
- Donald Olding Hebb. *The organization of behavior: A neuropsychological theory*. Psychology press, 2005.
- Geoffrey E. Hinton. Some demonstrations of the effects of structural descriptions in mental imagery. *Cogn. Sci.*, 3:231–250, 1979.
- Geoffrey E. Hinton. How to represent part-whole hierarchies in a neural network. *Neural Computation*, 35:413–452, 2021.
- Geoffrey E Hinton, Sara Sabour, and Nicholas Frosst. Matrix capsules with em routing. In *International conference on learning representations*, 2018.
- Eugene M. Izhikevich. Polychronization: Computation with spikes. *Neural Computation*, 18: 245–282, 2006.
- Jon H Kaas. Topographic maps are fundamental to sensory processing. *Brain Research Bulletin*, 44(2):107–112, 1997. ISSN 0361-9230. doi: [https://doi.org/10.1016/S0361-9230\(97\)00094-4](https://doi.org/10.1016/S0361-9230(97)00094-4). URL <https://www.sciencedirect.com/science/article/pii/S0361923097000944>.
- Peter König, Andreas K Engel, and Wolf Singer. Integrator or coincidence detector? the role of the cortical neuron revisited. *Trends in neurosciences*, 19(4):130–137, 1996.
- L. Kozachkov, Jean-Jacques E. Slotine, and Dmitry Krotov. Neuron-astrocyte associative memory. 2023. URL <https://api.semanticscholar.org/CorpusID:265157552>.
- Charles C Lee and S Murray Sherman. Drivers and modulators in the central auditory pathways. *Frontiers in neuroscience*, 4:14, 2010.
- T. S. Lee and Mai Lin Nguyen. Dynamics of subjective contour formation in the early visual cortex. *Proceedings of the National Academy of Sciences of the United States of America*, 98 4:1907–11, 2001.
- Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. *Advances in Neural Information Processing Systems*, 33:11525–11538, 2020.
- Sindy Löwe, Phillip Lippe, Maja Rudolph, and Max Welling. Complex-valued autoencoders for object discovery. *arXiv preprint arXiv:2204.02075*, 2022.
- Sindy Lowe, Phillip Lippe, Francesco Locatello, and Max Welling. Rotating features for object discovery. *ArXiv*, abs/2306.00600, 2023. URL <https://api.semanticscholar.org/CorpusID:258999615>.
- Keyvan Mahjoory, Jan-Mathijs Schoffelen, Anne Keitel, and Joachim Gross. The frequency gradient of human resting-state brain oscillations follows cortical hierarchies. *eLife*, 9, 2019.
- Christoph von der Malsburg. The correlation theory of brain function. In *Models of neural networks*, pp. 95–119. Springer, 1994.
- Nikola T. Markov, Julien Vezoli, Pascal J. P. Chameau, Arnaud Y. Falchier, René Quilodran, Cyril Huisoud, Camille Lamy, Pierre Misery, Pascale Giroud, Shimon Ullman, Pascal Barone, Colette Dehay, Kenneth Knoblauch, and Henry Kennedy. Anatomy of hierarchy: Feedforward and feedback pathways in macaque visual cortex. *The Journal of Comparative Neurology*, 522:225 – 259, 2013.

- Luca Mazzucato, Alfredo Fontanini, and Giancarlo La Camera. Dynamics of multistable states during ongoing and evoked cortical activity. *The Journal of Neuroscience*, 35:8214 – 8231, 2015.
- Jan Melchior and Laurenz Wiskott. Hebbian-descent. *ArXiv*, abs/1905.10585, 2019.
- Christoph Miehl, Sebastian Onasch, Dylan Festa, and Julijana Gjorgjieva. Formation and computational implications of assemblies in neural circuits. *The Journal of Physiology*, 601, 2022. URL <https://api.semanticscholar.org/CorpusID:252109958>.
- Wilten Nicola and Claudia Clopath. A diversity of interneurons and hebbian plasticity facilitate rapid compressible learning in the hippocampus. *Nature Neuroscience*, 22:1168 – 1181, 2019.
- Georg Northoff and Zirui Huang. How do the brain’s time and space mediate consciousness and its different dimensions? temporo-spatial theory of consciousness (ttc). *Neuroscience & Biobehavioral Reviews*, 80:630–645, 2017.
- Georg Northoff and Federico Zilio. Temporo-spatial theory of consciousness (ttc) – bridging the gap of neuronal activity and phenomenal states. *Behavioural Brain Research*, 424, 2022.
- Christos H. Papadimitriou, Santosh S. Vempala, Daniel Mitropolsky, Michael Collins, and Wolfgang Maass. Brain computation by assemblies of neurons. *Proceedings of the National Academy of Sciences of the United States of America*, 117:14464 – 14472, 2019.
- Jing Pei, Lei Deng, Sen Song, Mingguo Zhao, Youhui Zhang, Shuang Wu, Guanrui Wang, Zhe Zou, Zhenzhi Wu, Wei He, Feng Chen, Ning Deng, Si Wu, Yu Wang, Yujie Wu, Zheyu Yang, Cheng Ma, Guoqi Li, Wentao Han, Huanglong Li, Huaqiang Wu, Rong Zhao, Yuan Xie, and Luping Shi. Towards artificial general intelligence with hybrid tianjic chip architecture. *Nature*, 572:106 – 111, 2019.
- Joshua C. Peterson, David D. Bourgin, Mayank Agrawal, Daniel Reichman, and Thomas L. Griffiths. Using large-scale experiments and machine learning to discover theories of human decision-making. *Science*, 372:1209 – 1214, 2021. URL <https://api.semanticscholar.org/CorpusID:235396315>.
- Hubert Ramsauer, Bernhard Schaf, Johannes Lehner, Philipp Seidl, Michael Widrich, Lukas Gruber, Markus Holzleitner, Milena Pavlović, Geir Kjetil Ferkingstad Sandve, Victor Greiff, David P. Kreil, Michael Kopp, Günter Klambauer, Johannes Brandstetter, and Sepp Hochreiter. Hopfield networks is all you need. *ArXiv*, abs/2008.02217, 2020. URL <https://api.semanticscholar.org/CorpusID:220968978>.
- Rajesh P. N. Rao and Dana H. Ballard. Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. *Nature Neuroscience*, 2:79–87, 1999.
- Rajesh P. N. Rao, Dimitrios C. Gklezakos, and Vishwas Sathish. Active predictive coding: A unified neural framework for learning hierarchical world models for perception and planning. *ArXiv*, abs/2210.13461, 2022. URL <https://api.semanticscholar.org/CorpusID:253107762>.
- Pieter R. Roelfsema. Solving the binding problem: Assemblies form when neurons enhance their firing rate—they don’t need to oscillate or synchronize. *Neuron*, 111:1003–1019, 2023. URL <https://api.semanticscholar.org/CorpusID:257956926>.
- Peter J Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20:53–65, 1987.
- S Murray Sherman and RW Guillery. On the actions that one nerve cell can have on another: distinguishing “drivers” from “modulators”. *Proceedings of the National Academy of Sciences*, 95 (12):7121–7126, 1998.
- Wolf Singer. Distributed processing and temporal codes in neuronal networks. *Cognitive Neurodynamics*, 3:189 – 196, 2009. URL <https://api.semanticscholar.org/CorpusID:13388189>.

- Wolf Singer. The cerebral cortex: A delay-coupled recurrent oscillator network? In *Reservoir Computing*, 2021.
- Nelson Spruston. Pyramidal neurons: dendritic structure and synaptic integration. *Nature Reviews Neuroscience*, 9(3):206–221, 2008.
- Andrew Stockman. Cone fundamentals and cie standards. *Current Opinion in Behavioral Sciences*, 30:87–93, 2019.
- Shuyang Sun, Xiaoyu Yue, Song Bai, and Philip H. S. Torr. Visual parser: Representing part-whole hierarchies with transformers. *ArXiv*, abs/2107.05790, 2021.
- Catherine Tallon-Baudry and Olivier Bertrand. Oscillatory gamma activity in humans and its role in object representation. *Trends in cognitive sciences*, 3(4):151–162, 1999.
- Shiming Tang, Yimeng Zhang, Zhihao Li, Ming Li, Fang Liu, Hongfei Jiang, and Tai Sing Lee. Large-scale two-photon imaging revealed super-sparse population codes in the v1 superficial layer of awake monkeys. *eLife*, 7, 2018.
- Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Neural Information Processing Systems*, 2017. URL <https://api.semanticscholar.org/CorpusID:13756489>.
- Jonathan D Victor and Keith P Purpura. Nature and precision of temporal coding in visual cortex: a metric-space analysis. *Journal of neurophysiology*, 76(2):1310–1326, 1996.
- Christoph Von der Malsburg. The what and why of binding: the modeler’s perspective. *Neuron*, 24(1):95–104, 1999.
- Johan Wagemans, Jacob Feldman, Sergei Gepshtein, Ruth Kimchi, James R. Pomerantz, Peter A van der Helm, and Cees van Leeuwen. A century of gestalt psychology in visual perception: II. conceptual and theoretical foundations. *Psychological bulletin*, 138 6:1218–52, 2012.
- Xiaoqin Wang. Cortical coding of auditory features. *Annual review of neuroscience*, 41:527–552, 2018.
- Jeremy M. Wolfe. Forty years after feature integration theory: An introduction to the special issue in honor of the contributions of anne treisman. *Attention, Perception, & Psychophysics*, 82:1–6, 2020.
- Jeremy M. Wolfe. Guided search 6.0: An updated model of visual search. *Psychonomic Bulletin & Review*, 28:1060 – 1092, 2021.
- Yujie Wu, Lei Deng, Guoqi Li, Jun Zhu, and Luping Shi. Spatio-temporal backpropagation for training high-performance spiking neural networks. *Frontiers in Neuroscience*, 12, 2018.
- Danfei Xu, Yuke Zhu, Christopher B Choy, and Li Fei-Fei. Scene graph generation by iterative message passing. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5410–5419, 2017.
- Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications, 2023.
- Hao Zheng, Hui Lin, Rong Zhao, and Luping Shi. Dance of snn and ann: Solving binding problem by combining spike timing and reconstructive attention. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 31430–31443. Curran Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/cba76ef96c4cd625631ab4d33285b045-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/cba76ef96c4cd625631ab4d33285b045-Paper-Conference.pdf).
- Hao Zheng, Hui Lin, and Rong Zhao. GUST: Combinatorial generalization by unsupervised grouping with neuronal coherence. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=9005cvFzkZ>.

## A APPENDIX

In the main text, we leave out several details to make things self-contained within limited pages. These details are organized in the Appendix as followings:

### A.0.1 GENERAL DISCUSSIONS

We first discuss the **“Limitation”**, Broader Impact, and potential Future Work of this paper in Section.A.1 and Section.A.2.

### A.0.2 METRICS

Developing metrics to evaluate hierarchical relationship is a major challenge for this work. In the main text, we mainly provide the intuitions of what is the metric and what each score measures. In Section.A.3 (**“Metric”**), we provide a step-by-step introduction to how metrics are developed and how can these metrics be generalized to account for broader works in the future.

### A.0.3 HOW THE COMPOSER WORKS

While the intuition of the mechanism and the detailed architecture are introduced in the main text, several additional information might be helpful to further capture the picture of how the Composer works: (1) how the Composer is initialized? (2) What is the detailed configuration of time scale parameters and the robustness of the Composer to varied time scale settings (ablation study and sensitivity test). What is the contribution of time scale parameters to the Composer? (3) The training details of DAE and the contribution of DAE to the Composer. (4) To what extent the Composer is bio-plausible?

To resolve these potential questions, we show model details in Section.A.4 ~ Section.A.6, including **“Initialization”** in Section.A.4.1 and the **“Time Scale”** in Section.A.4.2. After that, we provide **“Training”** details in Section.A.5. Additional ablation studies of times scale parameters and DAE are provided in Section.A.4.3~Section.A.4.4, and Section.A.5.5 respectively. Then we also enumerate all the **“Biological Motivations”** and neural correlates of the Composer in Section.A.6.

### A.0.4 ADDITIONAL DISCUSSION OF THE MAIN-TEXT EXPERIMENTS AND ADDITIONAL RESULTS

If the reader has confusions on the detailed setting about the main text experimental figures. We provide supplementary discussions in Section.A.7.

### A.0.5 MORE VISUALIZATIONS

We lastly illustrate more detailed visualization results in Section.A.7.3 on different datasets (Fig.19 to Fig.26). Failure cases are included.

### A.0.6 ANIMATION VIDEO

Since the dynamical process is best illustrated by dynamical visualizations, we provide a zip file containing videos visualizing the dataflow of the Composer in SI , about 60MB to be downloaded.

### A.1 LIMITATION AND FUTURE WORK

In this section, we highlight several limitations that could be addressed in future works.

**Pixel-wise SCS as low-level feature layer and the Scalability.** In this work, the part-whole hierarchy is represented and evaluated based on the pixel-wise spike coding space (SCS), which has a one-to-one mapping to the image’s pixel space (topographical mapping). Shared pixel-wise spatial map of  $s_1, s_2, GT_1, GT_2, x, \gamma_1, \gamma_2, \Gamma$  to the input image  $x$  makes the spiking representation and the top-down attention much more interpretable so that the part-whole hierarchy could be explicitly evaluated and visualized as Fig.6, Fig.7. Also different object-level spiking pattern can be compared to compute the NEST Score. While topographical mapping (a one-to-one spatial relation between internal representation to the external physical world) is a common feature of the cortical representation, the representation (at each spatial point) could be more abstract. For example, instead of a binary neuron associated with the ‘presence’ of an object occupying the location, a population of binary neurons can be assigned to each location, so that different aspects of features (associated with the object at the corresponding location) could be accounted for by the population vector. Assigning a population of neurons with locations is similar to the capsule idea in Capsule Network Hinton et al. (2018) and the mini-column organization in GLOM Hinton (2021). This suggests a future direction to combine the neuronal synchrony (temporal coherence) with identical islands of neurons (spatial coherence) in the GLOM architecture.

Besides, the object could be represented in the latent layer of the DAE instead of the pixel-level SCS in each column, similar to Locatello et al. (2020). Representing objects in the latent layer could enable transforms between levels as in GLOM, by replacing pixel-wise gating in this paper into a neural network that potentially parameterizes a coordinate transformation. It is notable that the underlining hypothesis is that all feature-levels and object-levels share a “topographical” spatial map. The difference is what is being represented at each location (along feature level) and how representations at different locations are correlated (along object-level). This is a natural generalization of the representation hypothesis in Fig.1 in the Main Text.

Besides, the input to the SCS of the lowest level is not restricted to be the pixel-level image but could be the output feature map of an encoder, which is called tokenization in Agglomerator Garau et al. (2022). This kind of generalization has also been discussed in Hinton (2021). Notably, Lowe et al. (2023) scales the original “toy model” to represent objects in real-world dataset by applying the Visual Transformer (ViT) as a front-end encoder. By applying the ViT, original model only need to deal with much lower-dimensional representations, which have more similar structure as synthetic datasets. Lowe et al. (2023) provides the insights that seemingly toy models can be scaled to account for real-world dataset by applying proper front-end encoders. For Composer, a general-front end is also applicable and how to scale the model to account for neural syntax of more complicated scenes is a line of promising future works.

In sum, the limitation of the part-whole hierarchy as a pixel-level relationship in pixel-level SCS could potentially be generalized in three directions: (1) To allocate each location a column of spiking neurons to form a “representation column” at each location. (2) replacing the simplified cross-level interaction between pixel-wise SCSs as a proper neural network between latent layers. (3) the input to the model could be generalized to the tokenized embeddings from the front-end encoder. Notably, all these generalizations are compatible with the Composer and could be explored as future works.

**Coordination transformation.** As originally motivated in Hinton (1979) and restated in Hinton (2021), part-whole hierarchy contains two challenges: (1) the dynamical emergence of the part-whole tree structure and (2) the implementation of a part-whole coordinate transformation. The insight behind this paper is that the first challenge is the core challenge of the problem while the latter one could be solved by implementing the transformation as a neural network. In other words, the flexible forming of a symbolic tree structure (capable of capturing the basic nested part-whole relationship) within a pure neural network is the hard problem that challenges the neural network models. The second problem, implementing a coordinate transformation, is more compatible with the neural networks: such transformation could be realized as a (feedforward) neural network.

In this work, we focus on how to represent the part-whole hierarchy within a pure neural network model through emergent nested neuronal coherence. The coordinate of objects is assumed to be already aligned in a pixel-wise manner between the whole-level and part-level so that the coordination transformation is reduced to the inclusion mapping (similar to identity mapping). However, since the

parsing tree is realized within a pure neural network, the mechanism is compatible with more general coordination transformation: it could be realized by replacing the pixel-wise gating with a trainable neural network, which parameterizes the coordinate transformation.

**How many level are there?** In the representation hypothesis, we separate the entire representation space into discrete number of object levels. Therefore, different levels are explicitly distinguished before-hand. However, different from “discrete slots to represent object”, the discretization of representation into levels is a reasonable setting. The underlining hypothesis is that while object representation can be diverse and uncertain, human-vision only accounts for very limited finite number of levels at each instant. Although there are many object levels in the external world (from galaxy to atom), only a limited fraction of these levels are mapped to the internal representation space to form perception. For example, Hinton (2021) argued that 5 level is enough to account for human vision. The part-level and whole-level is a relative concept, so that a given object can be represented in either level based on the context. In other word, the internal finite number of object levels can be flexibly reused to account for different external object levels. This argument is the motivation to discretize the representation space explicitly into object-levels. But at each object level, objects are represented in a distributed manner, by neuronal coherence, instead of slots. As a result, the hierarchical relationship is also represented in a distributed manner.

In this paper, we show the part-whole hierarchy of two levels: whole and part. However, this minimal structure could be naturally extended to account for more levels, since the form of interaction between levels is mostly irrelevant to how many levels are there or which level it is in. All levels could share a unified form of cross-level interaction and within-level interaction. Therefore, by stacking the columns along the hierarchy, more levels are accounted for. As discussed in Hinton (2021), up to five levels are sufficient to realize human-like vision.

**Synthetic image.** In this paper, we use synthetic images to demonstrate how to represent the part-whole hierarchy. The benefit of using the synthetic image is that: (1) a common sense reasonable part-whole relationship is known beforehand as ground truth, therefore it is more convenient to explicitly evaluate the representation and test the capability. (2) The ground truth assignment of objects (part/whole) is known, which could be utilized to evaluate the neuronal coherence. The weak side of a ground truth is that such explicit assignment of part-whole ignores the ambiguity of parsing the scene: the parsing could depend on many factors like prior knowledge, attention, goal, internal state, and so on. Besides, parsing a real-world image without explicit part-whole hierarchy might be challenging for other reasons (overlap, background, etc.). However, recent models Hinton et al. (2018); Sun et al. (2021); Garau et al. (2022) that claim to solve the part-whole problem actually resemble performing hierarchical feature extractions. Such confusion is partly due to the ambiguity of the part-whole relation, object definition and the complexity of features in the real-world images, which confuse the symbolic structure. Therefore, taking the present status of the problem<sup>1</sup> and the challenges the problem implicates<sup>2</sup> into account, one desirable roadmap is to focus on explicit evaluation based on synthetic data first (so that it is easier to interpret whether the mechanism works) and then gradually generalize the outcome to increasingly complex datasets in the future.

**Learning scheme.** In this work, we treat the ‘sense’ of what the object should look like as prior knowledge embedding in the parameters of DAE’s weight, which in turn determines the dynamical property of the Composer. Indeed, such prior of prototype is needed for humans to parse a visual scene as well. For example, given a visual scene of a house (Fig.1a in the main text), a human observer should have already had the concept of the door, the window, and the roof in their mind, so that a house is parsed in the way Fig1.b (main text) shows. Therefore, in this work, we consider how a parsing structure could emerge as hierarchical neuronal coherence in a pure neural network given the prior knowledge of objects. On the one hand, some of the priors are indeed hard-wired in the brain through a long period of evolution (related to Gestalt psychology Wagemans et al. (2012), like proximity, similarity, enclosure, continuation, closure, symmetry, common fate, etc.); on the other hand, some of the others may be gradually learned during evolution. Therefore, the learning scheme could be updated to capture how the part-whole hierarchy could emerge during the unsupervised perception of the multi-object world. One recent work on DASBE shows that the

<sup>1</sup>Representing the part-whole hierarchy in a pure neural network is still an unsolved problem Hinton (2021)

<sup>2</sup>In essence, the part-whole problem requires a general solution to the sub-neuro-symbolic architecture and to realize hierarchical split of computational problem in a divide-and-conquer way. This paper explores the temporal aspect of the solution.



pre-training Zheng et al. (2022) in DASBE can be naturally generalized to end-to-end unsupervised learning Zheng et al. (2023) to predict the external input. The oscillatory cell assemblies also emerges, without supervision or any prior knowledge! The insight is that: the model architecture in this paper (delay-coupled column organization, time scale hierarchy) could be regarded as inductive bias to encourage a hierarchically factorized representation during the unsupervised training of the whole model as a recurrent spiking neural network to reconstruct what it sees on average during a temporal period. More details are discussed in Section.A.5.4. It is promising future work to figure out whether the “neural syntax” can be learned by simply predicting the visual input. This picture is consistent with the neuroscientific hypothesis of “Cognition from Action”Buzsáki et al. (2014). Therefore, we could treat the Composer as a constructed solution of the “neural syntax” by combining machine learning and neuroscience, and this solution provides the insights on how to learn the solution from data via unsupervised training, similar to from constructing the “neural words” in Zheng et al. (2022) to learning the “neural words” in Zheng et al. (2023). Taken together, we hope the Composer provide a new paradigm to combine machine learning models<sup>3</sup> and neuroscientific hypothesis (cell assemblies and cortical dynamics) to (1) solve part-whole hierarchy in machine learning field (in a bio-plausible way) and (2) understand how neural syntax emerges by predicting the world (action or attention).

## A.2 BROADER IMPACT

On the positive side, the model parses objects with neuronal coherence in the spike coding space composed of spiking neurons without explicit supervision. The mechanism by principle is not limited to a certain modality or certain object type. Thus, it may help develop human-like perception systems. Besides, with biological relevant features (eg. delayed coupling) and phenomena (eg. synchrony / neuronal coherence), the model may also act as a data-driven biological model to understand the perception process in the brain.

On the negative side, since the model is not supervised, it is harder to control what it learns. The current model is only trained on simple synthetic datasets and learns to group at the superficial pixel level, therefore the representation is highly explainable. However, grouping in latent space on real-world datasets requires to develop evaluation and visualization methods to make the representation in latent space more understandable. We believe this may serve as a step toward more transparent and interpretable predictions.

---

<sup>3</sup>DAE can be generalized into more advanced models like BERT, Diffusion model, etc, which also perform denoising.

### A.3 METRIC

In the main text, we introduced that the cell assemblies are treated as clusters of spike trains and therefore, the neuronal coherence is measured as the inner-cluster coherence of clusters, based on ground truth assignments. In this section, we further unfold how the metrics for quantitative evaluations are defined based on the Silhouette Score, including the Part Score, Whole Score, and Nest Score. Since the Silhouette Score is based on the similarity measure among samples, we first introduce how the similarity among spike trains is measured, where the Victor-Purpura metric shows up (Section.A.3.1). Then we introduce the Silhouette Score and how it could be extended to account for varied aspects of the part-whole representation (Section.A.3.2~Section.A.3.4). Detailed perturbation study of metrics is included in Section.A.3.5. Finally, we discuss how the metrics can be generalized to evaluate other models with similar attempts to group distributed representation into objects by certain similarity measures (identical islands of vectors in GLOM Hinton (2021)), in Section.A.3.7. Hyper-parameters for evaluation is shown in Table.1.

#### A.3.1 HOW TO MEASURE THE DISTANCE BETWEEN SPIKE TRAINS: VICTOR-PURPURA METRIC

The Victor-Purpura metric (VP-metric) is a classical non-Euclidean metric to measure the distance between arbitrary spike trains for evaluating the temporal coding in the visual cortex Victor & Purpura (1996). The motivation is that spike trains can have varied length and of binary value, so that the distance measure (like Euclidean metric) defined on fixed-length real-valued vector does not apply, especially to capture the precise temporal structure of spikes. The idea is that spike train can be treated as (binary) “Strings” and the distance from one spike train to the other can be defined as the number (cost) of “operations” to transform one spike train to the other (“Edit distance” of Strings).

In Victor & Purpura (1996), three types of operations are identified (Fig.9): 1. add a spike (cost=1); 2. delete a spike (cost=1); 3. shift a spike for length  $\Delta t$  (cost= $\Delta t/\tau$ ), where  $\tau$  is a parameter to control the temporal precision of the spiking code (or the temporal sensitivity of the metric). By sequentially applying the three operations ( $T(u)$  in Fig.9), a spike train can be transformed to the other. The Victor-Purpura distance is defined as the minimal cost to transform a spike train to the other (Fig.9):

$$D_{VP}(s_i, s_j; \tau^{-1}) = \min_T \left( \sum_{u=1}^{|T|} \text{cost}(T(u)) \right) \quad (8)$$

$$T(u) \in \{delete, add, shift\} \quad (9)$$

$$\text{cost}(delete) = \text{cost}(add) = 1 \quad (10)$$

$$\text{cost}(shift, \Delta t; \tau) = \Delta t/\tau \quad (11)$$

where  $T(u)$ ,  $u = 1 \dots |T|$  is a sequence of basic transformations to transform  $s_i$  to  $s_j$  (or vice versa). The costs of the three basic transformations are different. The most special one is the shift operation: there is a time scale parameter to control the punishment of shifting the spike to its neighbourhood.

Notably, it is proved that the definition satisfies the three principles of a metric: positivity, symmetry, and triangle inequality Victor & Purpura (1996). Thus, it induces a metric space of arbitrary spike trains, even if not embedded in a vector space of specified dimension. Since spike trains are non-Euclidean in nature, the VP-metric provides a more direct measure of these entities. The minimal cost is computed through a dynamic programming method.

A desirable feature of VP-metric is that the parameter  $\tau$  explicitly controls the temporal sensitivity of the metric and the expected temporal precision to be considered. If the  $\tau$  is chosen to be  $\infty$ , then shifting a spike will cause no cost ( $1/\infty \sim 0$ ). Thus the distance is exclusively due to spike count<sup>4</sup>, therefore spiking rate is measured. If the  $\tau$  is chosen to be 0, then it measures the number of spikes

<sup>4</sup>All transforming cost comes from adding/deleting spikes

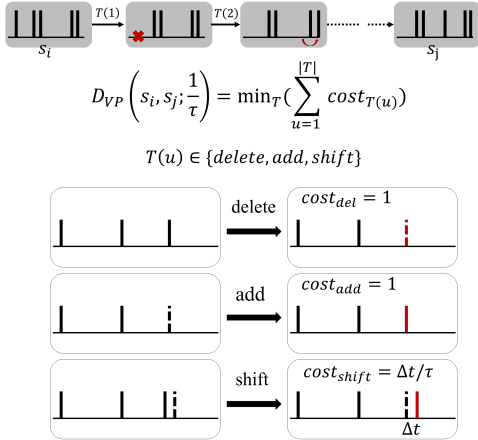


Figure 9: The Victor-Purpura metric. Top: a sequence of transformations; Middle: the distance between  $s_i$  and  $s_j$  is defined as minimum cost to transform one to the other; Bottom: three basic operations and their cost. Dashed bars are imaginary spikes which has been removed or has not been added. Red bars highlight the results due to the transformation.

that are not in absolute synchrony<sup>5</sup>. So small  $\tau$  measures the spike train distance based on the very precise temporal synchrony structure. By varying the  $\tau$ , it is plausible to find the optimal coding scheme of the visual cortex Victor & Purpura (1996). In sum,  $\tau$  is treated as a timescale parameter, to control the precision of temporal coding. In this paper, the part level is evaluated with smaller  $\tau$  (part-level time scale) while whole level are evaluated with larger  $\tau$  (whole-level time scale). The nestedness between the levels is evaluated with the time scale of the whole level (the child node should stay within the time window of their parents.)

### A.3.2 GENERAL COHERENCE MEASURE OF CLUSTERS: SILHOUETTE SCORE

The Silhouette coefficient Rousseeuw (1987) is a score to evaluate the quality of clustering by measuring the inner-cluster coherence. Given a clustering assignment, the score is calculated using average intra-cluster distance (a) and average nearest-cluster distance (b). The score is computed as  $(b - a) / \max(a, b)$ . The document can be found at [https://scikit-learn.org/stable/modules/generated/sklearn.metrics.silhouette\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.silhouette_score.html). The best score is 1 and the worst score is -1. The values near 0 indicate overlapping clusters.

### A.3.3 VICTOR-PURPURA METRIC + SILHOUETTE SCORE

How could the neuronal coherence be measured? Given ground truth assignments of clusters (each neuron belongs to which part or whole object), neuronal coherence (synchrony) is measured as the inner-cluster coherence of spike trains: the “inner-cluster” similarity and “inter-cluster” separability.

Specifically, if we take (1) each neuron as a **sample**, (2) the spike train of each neuron as **features**, (3) the ground truth assignment as the **clustering assignment** (each pixel on ground truth correspond to a neuron in SCS), (4) the VP-metric as the **distance measure**, then, the inner-cluster coherence (Silhouette Score) of the ground-truth-induced cluster assignment is exactly the coherence measure of neuronal coherence. In other words, the high Silhouette Score indicates that the spike trains of neurons of the same cell assembly are closer to each other in terms of VP-metric, which can be interpreted as neurons in the same cell assembly synchronizing better. Therefore, the VP-induced Silhouette Score sufficiently measures the grouping quality. The VP-induced Silhouette Score is also from -1 to 1. The best value is 1 (perfect grouping) and values near 0 indicate overlapping clusters (purely random firing without any temporal structure). Negative values generally indicate

<sup>5</sup>only total synchronous spike trains have 0 distance while the slight shift of spikes has cost 1

that a sample has been assigned to the wrong cluster, as a different cluster is more similar (neurons systematically synchronize to incorrect groups).

During the whole simulation, only a segment of simulation (after the convergence) is used for evaluating the Silhouette Score, See Table.1. The length of segment can be flexibly selected as long as it covers at least one oscillation period. In this work, we simply select one large enough length value.

#### A.3.4 SCORES TO MEASURE THE PART-WHOLE HIERARCHY

Coherence scores are all defined based on the VP-Silhouette Score and Ground Truth assignment.

**Part Score** is defined as the VP-Silhouette score with respect to the part-level spiking pattern and the part-level ground truth assignment:

$$\text{Part Score} = \text{Silhouette}(VP(sp_{k_1}, sp_{k_1}; \tau_p), \text{label}_1) \quad (12)$$

where  $sp_{k_1} \in \{0, 1\}^{(N, \tau)}$  means the total spike trains in level 1 (part-level).  $\tau$  is the length of each spike train for evaluation and  $N$  is the number of (activated) neurons at part level.  $VP(sp_{k_1}, sp_{k_1})$  is the distance matrix whose elements  $(i, j)$  are the VP-distance between the  $i$ -th spike train and the  $j$ -th spike train in the part level; The  $\text{label}_1$  means the ground truth assignment of neurons in level 1 (part-level). *Silhouette* is the Silhouette Score.  $\tau_p$  is close to the (integration) time constant of part-level ( $\tau_1$ ) (Table.1), controlling the temporal sensitivity of VP-metric (eq.8). Therefore, the Part Score measures the coherence level of the part level exclusively, independent of the activity in the whole level. Part Score indicates the quality of the grouping of tree nodes in the part level (Fig.12). For example, *Part - Score* = 1 indicates that neurons are synchronized perfectly into separated groups corresponding to the part objects. On the other hand, Part Score = 0 indicates that the neurons fire randomly and no temporal structure emerges (Fig.11~Fig.13, bottom). In rare cases, Part Score < 0 indicates that (on average) neurons with different assignments are synchronized and neurons with the same assignments are not synchronized (Fig.14, bottom). Equation.12 is equivalent to the definition of Part Score in the main text, but eq.12 reveals how the Part Score is practically implemented. Same for other cases.

**Whole Score** are similarly defined as:

$$\text{Whole Score} = \text{Silhouette}(VP(sp_{k_2}, sp_{k_2}; \tau_w), \text{label}_2) \quad (13)$$

where  $sp_{k_2} \in \{0, 1\}^{(N, \tau)}$  means the total spike trains of level 2 (whole-level).  $VP(sp_{k_2}, sp_{k_2})$  is the distance matrix whose elements  $(i, j)$  are the VP-distance between the  $i$ -th spike train and the  $j$ -th spike train in the whole level; The  $\text{label}_2$  means the ground truth assignment of neurons in level 2 (whole-level).  $\tau_w$  is close to the time constant of the whole-level  $\tau_2$ :  $\tau_w > \tau_p$  (Table.1). The Whole Score measures the grouping quality of tree nodes at whole-level (Fig.13).

On the one hand, the forming of tree nodes is a necessary condition to form the entire tree, so Part Score and Whole Score are important measures of the representation. On the other hand, since the Part Score and Whole Score measures the grouping in part/whole level independently (Fig.5cd in the main text or Fig.12 ~ Fig.13), they do not reveal the correlation between levels. For example, the emergent cell assemblies can be arbitrarily permuted or translated (together) without affecting the scores (Fig.5e in the main text or Fig.14). Obviously, such arbitrary operations are serious enough to destroy a well-defined tree structure (Fig.14).

Therefore, we provide the additional score to capture the cross-level coordination: the Nest Score.

**Nest Score** is defined as the coherence score between part-level spiking patterns and whole-level spiking patterns based on whole level assignments:

$$\text{Nest Score} = (4/3) \cdot \text{Silhouette}(VP(sp_{k_1}, sp_{k_2}; \tau_w), \text{label}_2) \quad (14)$$

where  $sp_{k_1}, sp_{k_2} \in \{0, 1\}^{(N, \tau)}$  means the total spike trains of level 1 (part-level) and level 2 (whole level) respectively.  $VP(sp_{k_1}, sp_{k_2})$  is the distance matrix whose elements  $(i, j)$  are the VP-distance

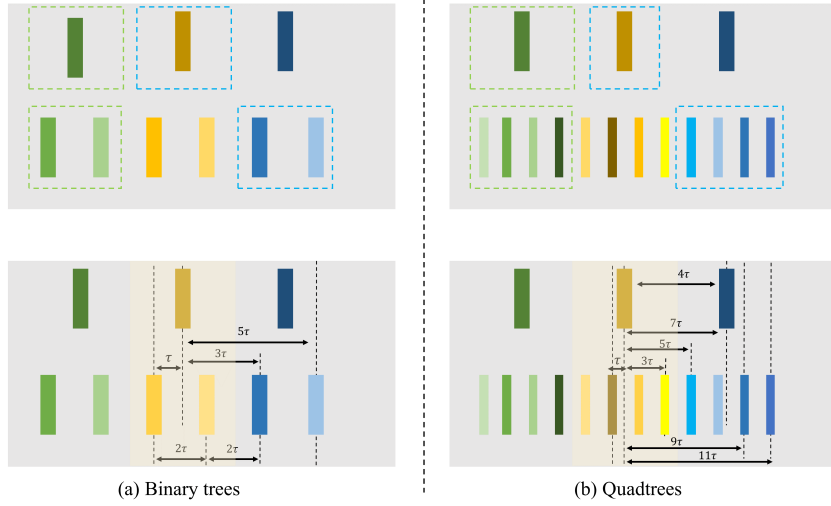


Figure 10: Illustration of the Nest Score. Each colored bar stands for a population of synchronized neurons (cell assemblies).

between  $i$ -th spike train in the part level and  $j$ -th spike train in the whole level.  $(4/3)$  is a normalization factor, introduced below.

Intuitively, the nestedness means that synchronized cell assemblies at the part-level are coordinated within the lifetime of whole-level cell assemblies (Fig.1 in the main text). Therefore, it could be formalized as the coherence measure between part-level spike trains and whole-level spike trains:  $(b - a) / \max(a, b)$ . Here,  $a$  is the average distance between part-level spike trains and whole-level spike trains that share the same whole-level assignment (Fig.10 top, green dashed box).  $b$  is the mean distance between a whole-level spike train and the nearest part-level cell assemblies that the spike train is not a part of (the averaged distance between sets, Fig.10 top, blue dashed box). This formulation is similar to the Silhouette, but replaces the  $VP(sp_{k_1}, sp_{k_1})$ ,  $VP(sp_{k_2}, sp_{k_2})$  to the  $VP(sp_{k_1}, sp_{k_2})$ . If the Nest Score is high, it means the part-level neurons  $s_1$  for the same whole-level objects ( $label_2$ ) is correlated to whole-level neurons  $s_2$ . This spatial-temporal structure can be translated into the nestedness of cell assemblies across levels.

However, since part-level and whole-level should not be exactly the same, the derived  $(b - a) / \max(a, b)$  do not reach the 1 in best cases. To normalize the score into the range of  $(-1, 1)$ , we compute a compensatory factor — the “magic number”:  $3/4$ . To simplify the problem, we assume that in the ideal parsing case, both part-level and whole-level spikes are synchronized perfectly and arranged uniformly along the time dimension (Fig.10 bottom). Assume the nearest time interval between part-level and whole-level spikes is  $\tau$  (Fig.10 bottom), then for a “perfect” binary tree (Fig.10 a):

$$a = (\tau + \tau) / 2 = \tau \quad (15)$$

$$b = (3\tau + 5\tau) / 2 = 4\tau \quad (16)$$

$$\text{Nest Score} = (b - a) / \max(a, b) = 3/4 \quad (17)$$

and for a “perfect” quadtree (Fig.10 b):

$$a = (\tau + \tau + 3\tau + 3\tau) / 4 = 2\tau \quad (18)$$

$$b = (5\tau + 7\tau + 9\tau + 11\tau) / 4 = 8\tau \quad (19)$$

$$\text{Nest Score} = (b - a) / \max(a, b) = 3/4 \quad (20)$$

Interestingly, in both ideal cases, the Nest Score is  $3/4$ . As a result, we normalize the derived Silhouette Score by a factor ( $3/4$ ), which is exactly the eq.14. Here binary tree accounts for SHOPS, Ts, and double-digit MNIST while quadtree accounts for the Squares. The validity of the normalization is confirmed in Fig.5 or Fig.14, where Nest Score achieves 1 in ideal cases.

Lastly, the reason why it is valid to use whole-level assignment to ground both part-level and whole-level neurons is that we have assumed a one-to-one spatial relation between part-level SCS and whole-level SCS (topographical mapping to the physical world), and the complete object has a 'copy' at each level along the hierarchy (main text). It is also a conventional assumption of the cortex Hinton (2021).

### A.3.5 PERTURBATION STUDY

In order to verify the proposed scores, we conduct a perturbation study in Fig.5 in the main text. Here, we provide more details on how the perturbation is made and more discussions about the experiment (Fig.11 ~ Fig.14).

Given an input image and its ground truth as in Fig.5a, we firstly 'artificially' build up the "perfectly" nested spike pattern in two oscillation periods (whole period ( $2 \cdot \tau_{total}$ ) is 54 time steps). More specifically, all cell assemblies are synchronized almost perfectly and arranged uniformly along the time axis as in Fig.10. Part-level cell assemblies are coordinated within the lifetime of whole-level cell assemblies (Fig.5a). It is the ideal case, with 0 perturbation level in Fig.5bcde or Fig.11 ~ Fig.14.

Then, for Fig.5bcd or Fig.11 ~ Fig.13, we randomly and independently perturb the timing of spikes into nearby time points:

$$t_i \longrightarrow t_i \pm \Delta t, \quad \Delta t \leq \tau \quad (21)$$

where  $t_i$  is the spike timing of  $i_{th}$  neuron,  $i \in N$ .  $N$  is the number of total considered neurons (part-level or whole-level or both levels).  $\tau$  is the timescale controlling the perturbation level. If  $\tau = 0$ , perturbation is zero. If  $\tau$  equals half the length of the oscillation period ( $0.5 \cdot \tau_{total}$ ), the perturbation will lead to pure random firings like Fig.11 ~ Fig.13, bottom. As a result, we define the perturbation level in Fig.5bcd as  $\tau / (0.5 \cdot \tau_{total})$ , ranging from 0% to 100%. A more detailed perturbation process is shown in Fig.11 to Fig.13. In Fig.5b or Fig.11, the perturbation is applied to both part and whole level so that all scores smoothly decrease from 1 to near 0. In Fig.5c or Fig.12, the perturbation is only applied to the part level, so that the Whole Score is not affected but both Part Score and Nest Score smoothly decrease from 1 to near 0. In Fig.5d or Fig.13, the perturbation is only applied to the whole level, so that the Part Score is not affected but both Whole Score and Nest Score smoothly decrease from 1 to near 0. In a word, in Fig.5c and Fig.5d (Fig.12 and Fig.13), we isolatedly verify the property of Part / Whole Score, which shows that they are capable of capturing the quality of node-level representation of a tree structure. In Fig.5b or Fig.11, we provide more common cases where both part and whole level degrades, which shows that three scores consistently measure the coherence of neuronal representation.

For Fig.5e or Fig.14, to isolatedly verify the role of Nest Score. We firstly build up perfect synchronized cell assemblies as in the perfect nested case (Fig.5a), but then perturb the timing of each 'cell assembly' at different levels. All spikes within the same synchronized cell assembly are perturbed with the same  $\Delta t$  and different cell assemblies are perturbed by independent  $\Delta t$  (Fig.5e top or Fig.14). Similarly, we define perturbation level as  $\tau / 0.5 \cdot (\tau_{total})$ , where  $\tau$  is the timescale controlling the perturbation. A more detailed perturbation process is shown in Fig.14. As shown in Fig.5e or Fig.14, the Nest Score decreases smoothly while Part / Whole Scores remain constant. Notably, the perturbation can lead to wrong hierarchical coordination: whole-level cell assemblies are synchronized with part-level cell assemblies of different assignments (incoherence). Thus, the Nest Score can decrease into values even lower than 0.

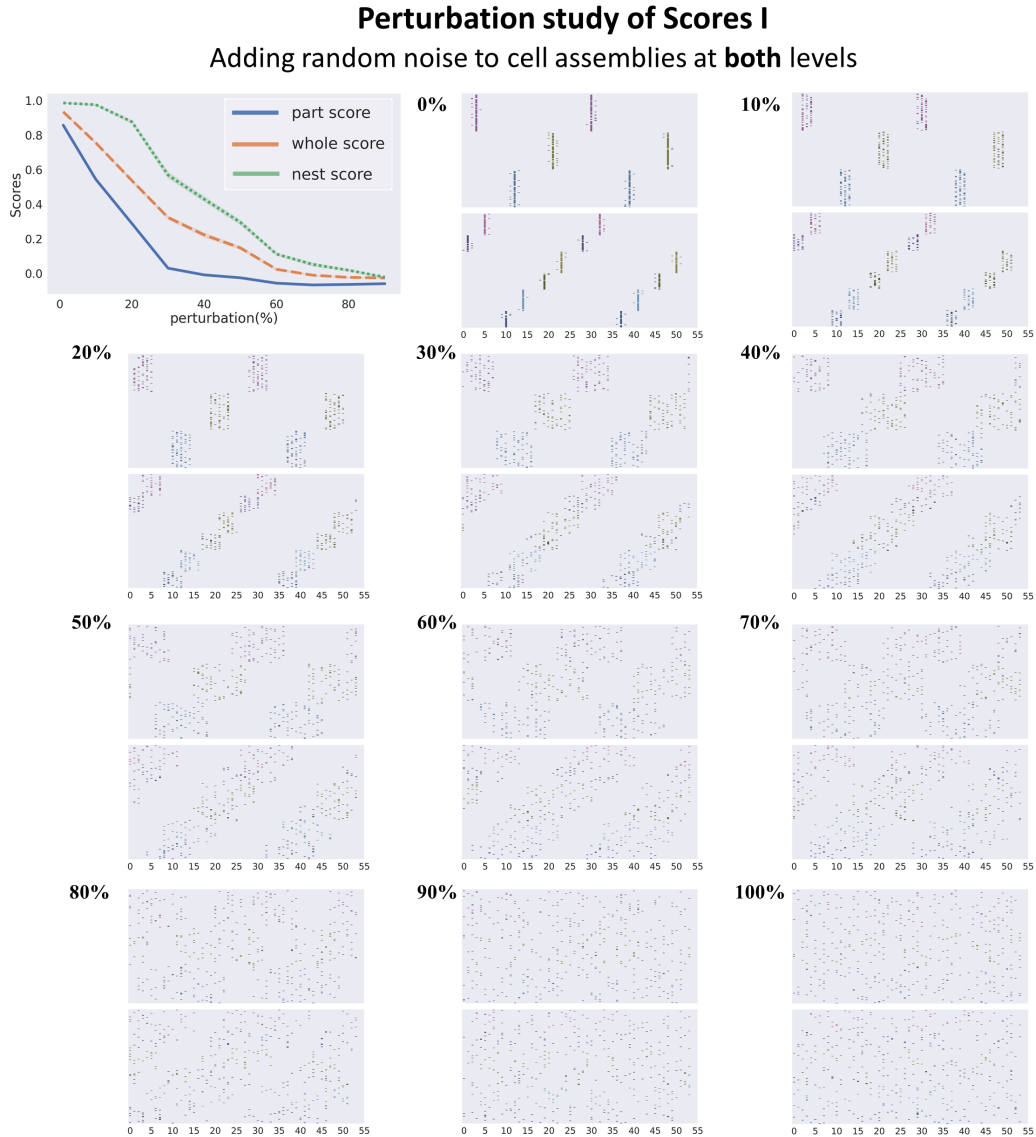


Figure 11: Perturbation study by adding increasing random noise to both part-level and whole-level spike patterns. Visualizations of the perturbed spiking pattern at different perturbation levels (0% to 100%) are shown, corresponding to the Fig.5b. Both part-level and whole-level gradually degrades into random firings.

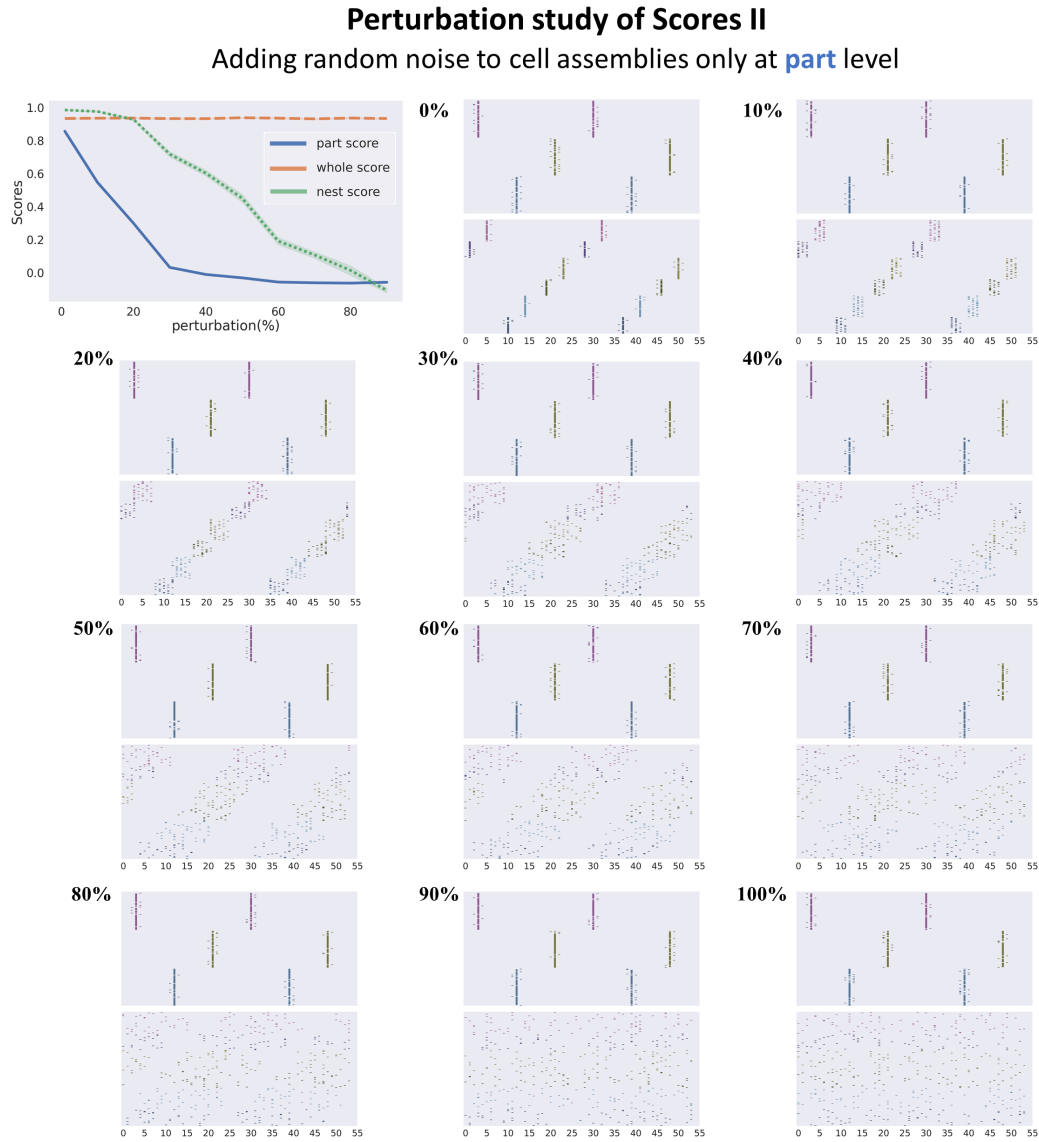


Figure 12: Perturbation study by adding increasing random noise to only part-level spike patterns. Visualizations of the perturbed spiking pattern at different perturbation levels (0% to 100%) are shown, corresponding to the Fig.5c. Part-level is degraded gradually while whole level remains unchanged.



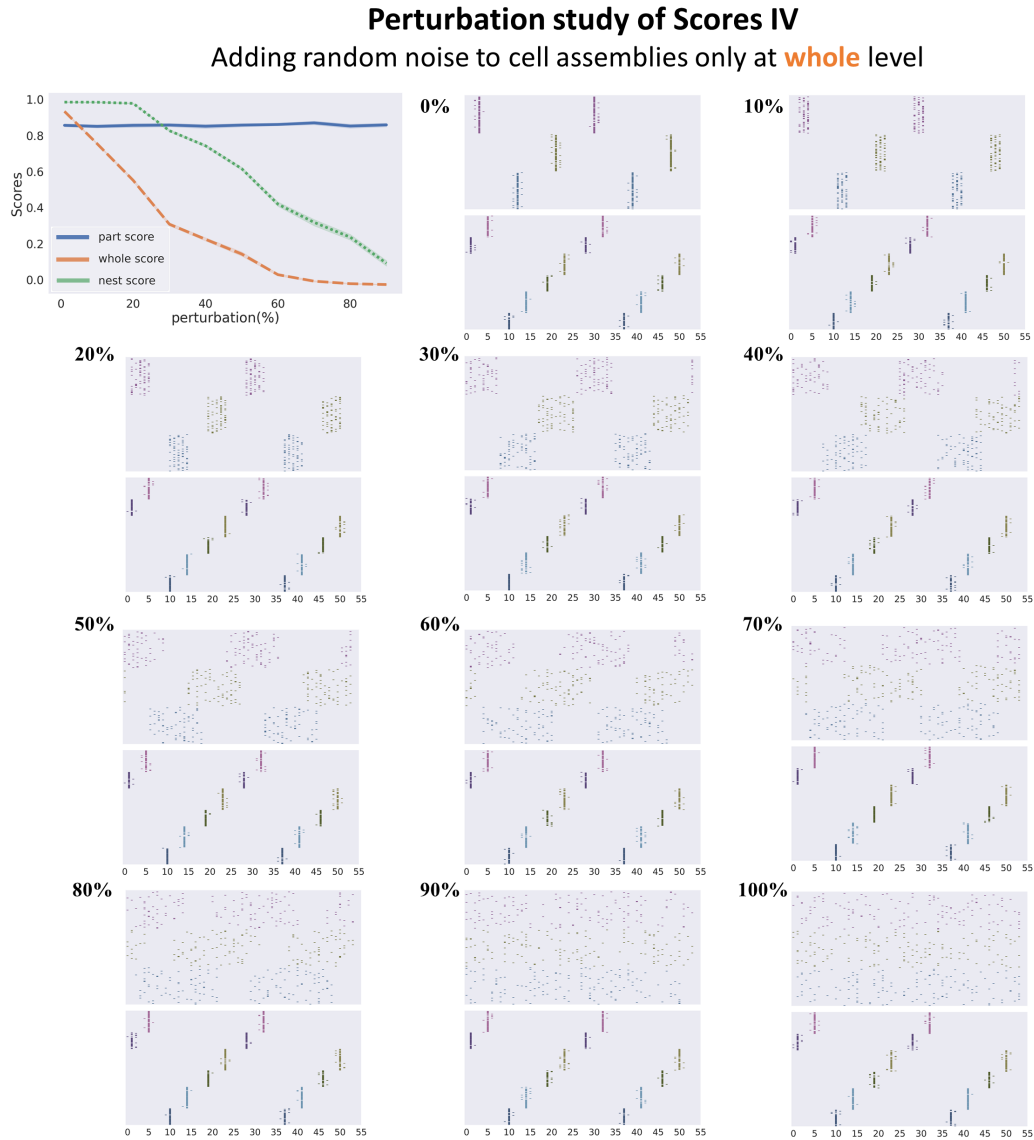


Figure 13: Perturbation study by adding increasing random noise to only whole-level spike patterns. Visualizations of the perturbed spiking pattern at different perturbation levels (0% to 100%) are shown, corresponding to the Fig.5d. Whole-level is degraded gradually while part level remains unchanged.

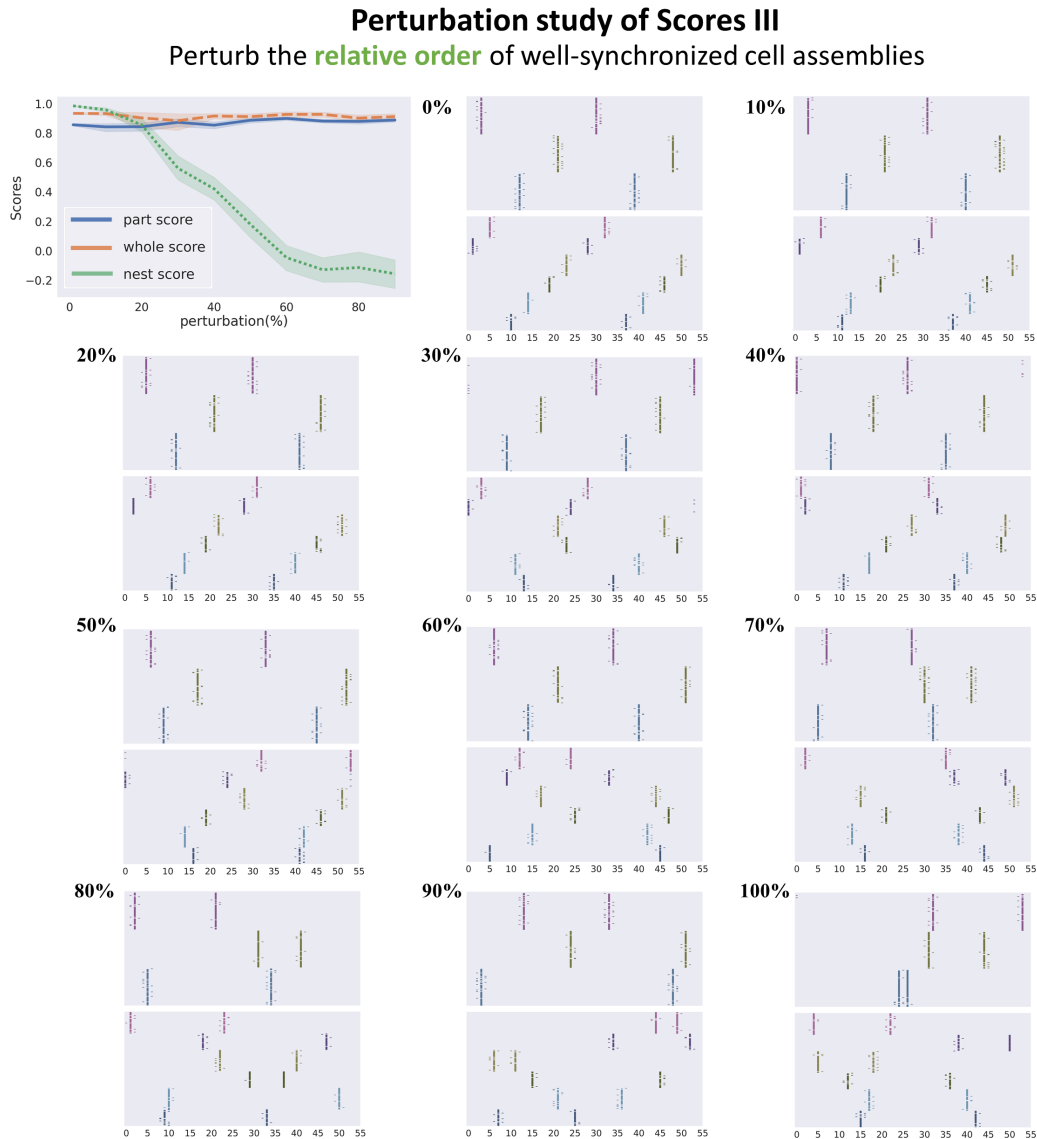


Figure 14: Perturbation study by adding increasing random noise to each assembly (both part-level and whole level). Visualizations of the perturbed spiking pattern at different perturbation levels (0% to 100%) are shown, corresponding to the Fig.5e. Relative coordination among cell assemblies is gradually changed while the synchronization of each cell assembly is unchanged.

Table 1: Parameters of the evaluation.  $\tau_p, \tau_w$  is the time scale parameter for computing the VP-distance in each case. Shift is the fixed modified time steps for computing Nest Score and for visualization. The segment length of spike trains for (1) computing scores ( $\tau_l$ ) and (2) for visualization is also shown.

Dataset	Ts	Squares	SHOPs	Double MNIST
$\tau_p$	2	2	3	2
$\tau_w$	6	7	6	4
shift	10	9	6	2
segment length (score, $\tau_l$ )	160	75	42	32
segment length (visualization)	200	100	100	70

### A.3.6 THE SHIFT: FROM SYNCHRONIZATION TO POLYCHRONIZATION

In the neural system of the brain, the synchronization matters since the coincident arrival of spike trains could have a much larger effect on the target neuron. Therefore, it is the “synchronization” in the viewpoint of the reader neuron that really matters Buzsáki (2010). However, due to the diverse axonal delay (tens of milliseconds) of different neurons, coincidentally arrived spikes are usually fired at different timings, yet with fixed temporal shifts. This phenomenon is called polychronization, or polychronized neuronal groups (PNG) Izhikevich (2006), which generalizes the concept of synchronization and is a more natural outcome of a real-world neural system, with potentially heterogeneous parameter settings. In other words, polychony and synchrony bear the same spirit of fixed temporal correlation, but polychronization could tolerate a fixed temporal shift among spike timings. While in an external observer’s viewpoint, two things are different, in the viewpoint of the internal readout neuron, both can be the same thing.

In other words, if we shift all timing patterns with a fixed shift parameter, it is equivalent to the original pattern in the sense that the temporal shift can be compensated by the fixed axonal delay when being read out by a downstream module. Motivated by this fact, such slight fixed shifts are compensated (ignored) before computing the Nest Score<sup>6</sup>. In other words, in the Composer, whole-level and part-level cell assemblies are allowed to have a slight fixed temporal shift (translation slightly along the time axis). The representation is regarded as unaffected as long as the shift is a constant (Table.1).

### A.3.7 GENERALIZE TO EVALUATE OTHER MODELS

The metric proposed in this paper is also applicable to neural models that exploit similarity or coherence measures to group neural representations into part-whole hierarchies. GLOM Hinton (2021) and GLOM-inspired Agglomerator Garau et al. (2022) is one interesting example. To measure the similarity among vectors, the Victor-Purpura metric is not needed anymore. Therefore, it is more direct to take each vector as a sample and the different dimensions of the vectors as the features in a clustering algorithm. In this way, three Silhouette-based coherence measures of spike trains could be naturally generalized to account for real-valued vectors. For example:

$$\text{Part Score} = \text{Silhouette}(D(l_p, l_p), \text{label}_1) \quad (22)$$

$$\text{Whole Score} = \text{Silhouette}(D(l_w, l_w), \text{label}_2) \quad (23)$$

$$\text{Nest Score} = \alpha \cdot \text{Silhouette}(D(l_p, l_w), \text{label}_2) \quad (24)$$

<sup>6</sup>Let  $s'_i$  denotes the shifted spiking patterns of the original spiking pattern  $s_i$ , then:  $s'_2(t) = s_2(t + \text{shift})$  while  $s'_1(t) = s_1(t)$ . In other words, we slightly shift whole-level assemblies “backwards” relative to part-level spiking patterns for computing the Nest Score and for visualization. In Fig.6, we shift the  $\gamma_i$  and  $\Gamma$  in the same way for the same reason. The shift time step is very small and is a constant (Table.1)

For GLOM Hinton (2021),  $l_p$  and  $l_w$  is the feature vector for each part-level column and whole-level column<sup>7</sup>. Besides, Euclidean metric for real-valued vector suffice for measuring the similarity:  $D$ .  $label_i$  is the ground-truth parsing<sup>8</sup>.  $\alpha$  is a normalization factor that can be determined based on the “ideal” case.

For complex-valued network Löwe et al. (2022), which could potentially form compositional structures by similarity among the phase of complex values<sup>9</sup>.  $l_p$  and  $l_w$  is the phase vector for part-level and whole-level. Similarly, Euclidean metric suffices for measuring the similarity among phase vectors:  $D$ .

Taken together, the proposed metric can be generalized to evaluated related models that (1) represent part-whole structure by certain type of coherence measure and (2) being tested on synthetic datasets where ground truth is available.

---

<sup>7</sup>An interesting point is that the islands of identical vectors in GLOM are parallel to the correlated cell assemblies in the Composer, as long as we take each temporally unfolded spike train as the (binary) vector in each GLOM’s column.

<sup>8</sup>We are aware that the proposed coherence metric still requires the ground truth ( $label_i$ ), which is a limitation to evaluate the performance on real-valued datasets. However, evaluating grouping or parsing without ground truth is a well-known big challenge for these fields and we leave this hard problem to future works.

<sup>9</sup>Currently, this line of work can not account for hierarchical object level since the grouping mechanism is realized as special activation function, which can not be flexibly modified to, for example, distinguish different object levels. In contrast, the time window of readout neurons in the Composer can be flexibly configured to shape different dynamics for different part-whole levels.

Table 2: Time scale constants of the Composer

Dataset	Ts	Squares	SHOPs	Double MNIST
$T$	3000	3000	3000	3000
$\tau_d$	80	75	42	16
$\tau_{\delta 1}$	36	36	20	16
$\tau_{\delta 2}$	35	35	20	15
$\tau_{r1}$	15	24	12	16
$\tau_{r2}$	14	24	12	15
$\tau_1$	2	2	3	2
$\tau_2$	6	12	6	8
$\tau_D$	18	30	10	8
$\tau_\Gamma$	15	16	8	8
$\tau_{d'}$	80	75	42	16

#### A.4 HOW THE COMPOSER WORKS

In the main text, we introduced the architecture of the Composer, the neuroscientific motivation of the architecture and the intuition of how the architecture generates the nested dynamics. However, while in the main text, we briefly mentioned that time scale parameters and the priors in DAE contribute to the “magic”, we leave out the details on parameter setting and training scheme into this Section and the following Section.

In this section, we provide more details about practical implementations on how the model works. We first clarify the initialization of the dynamics in Section A.4.1 and Table.3. Then, we discuss the time scale parameter settings (Table.2) and ablation of time scale parameters in Section.A.4.2 and Section.A.4.4. In next Section.A.5, we provide training details of the DAE and more ablations about DAE (Section.A.5.5). Taken together, a more complete picture of how different modules contribute to the Composer is unfolded.

To have a intuitive impression of the dataflows in the Composer, a zip file containing videos visualizing the dynamics of the Composer is provided in SI (60MB).

##### A.4.1 INITIALIZATION OF THE DYNAMICS

The dynamics of Composer is gated by the top-down feedback (output from DAE or integrated spikes from higher-level SCS). However, at the initial phase, what should be the value of DAE output or higher-level feedback? One simple solution is to initialize the feedback as (uniform) random noise.

To formulate, eq.1 and eq.4 in main text are slightly extended to clarify the initialization process:

$$\rho_1(t) = x \cdot (\gamma_1 \cdot \Gamma_1 + r_1(t) \cdot \epsilon_1) \quad (25)$$

$$\rho_2 = (\lambda \cdot x + (1 - \lambda) \cdot D) \cdot (\gamma_2 + r_2(t) \cdot \epsilon_2) \quad (26)$$

where, the term  $r_i(t) \cdot \epsilon_i$  is only for random initialization (See Table.3).  $\epsilon_i$  is sampled from uniform distribution  $U[0, 1]$  and  $r_i(t)$  is the temporary amplitude of the noise, which is decayed rapidly along the simulation (decay rate  $\sim 0.8$ , Table.3). In other words,

$$r_i(t) = r_i \cdot (0.8)^{-t/\tau_d}, i = 1, 2 \quad (27)$$

where  $r_i$  is the initial amplitude of the noise. During simulation, the noise is decayed every  $\tau_d$  time steps for simplicity.

##### A.4.2 TIME SCALES

The Composer is inspired by the neural syntax hypothesisBuzsáki (2010), which argues that the hierarchical organization of cell assemblies should be readout by hierarchical organization of time

Table 3: hyper-parameters of the Composer (besides time scale constants).  $g$  (eq.2) is the (inhibitory) gating effect of relative refractory period ( $\tau_\delta - \tau_r$ ), same for whole-level and part-level for simplicity.  $\lambda$  in eq.4 describes the skip connection.  $r_1, r_2$  describes the initialization (eq.27).

Dataset	Ts	Squares	SHOPs	Double MNIST
$g$	0.5	0.3	0.3	0
$\lambda$	0.3	0.4	0.4	0.4
noise decay	0.8	0.8	0.8	0.8
$r_1$	$\frac{1}{40}$	$\frac{2}{3}$	$\frac{1}{9}$	$\frac{1}{8}$
$r_2$	$\frac{1}{40}$	$\frac{2}{3}$	$\frac{1}{9}$	$\frac{1}{8}$

windows. In fact, it is likely in the Composer that the hierarchical organization of time window also encourages the emergence of hierarchical cell assemblies. Motivated by the hypothesis, spikes in SCS are all integrated within a narrow time window (implemented as a integration function  $I_i$  of certain time constant  $\tau_i$ ,  $i \in \{1, 2, D, \Gamma\}$ ) before the downstream processing. The time scale parameters related to the model are shown in Table.2, which has appeared in eq.1 to eq.7 in the main text.  $T$  is the entire simulation length.  $\tau_d$  is the coupling delay of the top-down feedback inside the column, shared for both part-level and whole-level.  $\tau_{\delta 1}$  is the total refractory period of part-level spiking neurons and  $\tau_{\delta 2}$  is that of the whole-level spiking neurons.  $\tau_{r1}$  is the absolute refractory period of part-level spiking neurons and  $\tau_{r2}$  is that of the whole-level spiking neurons.  $\tau_1$  is the integrative time window from part-level SCS to part-level DAE (eq.3) and  $\tau_2$  is that from the whole-level SCS to whole-level DAE (eq.6).  $\tau_D$  is the integrative time window from the part-level SCS to the whole-level SCS (eq.7, right).  $\tau_\Gamma$  is the time window of the top-down feedback from the whole-level SCS to the part-level SCS (eq.7, left).  $\tau_{d'}$  is the coupling delay of the cross-level top-down feedback from the whole-level SCS to the part-level SCS (eq.7, left). In this work, we set  $\tau_{d'} = \tau_d$  for simplicity. Roughly speaking, we have:

$$\tau_d = \tau_{d'} > \tau_{\delta 1} \sim \tau_{\delta 2} > \tau_{r1} \sim \tau_{r2} > \tau_2 \sim \tau_D \sim \tau_\Gamma > \tau_1 \quad (28)$$

Notably, part-level and whole-level columns are characterized by two timescale parameters (readout time window of SCS):  $\tau_1 < \tau_2$ , which softly shape the timescale (fast or slow) of the intra-column dynamics, which is inspired by the timescale hierarchy along the cortical hierarchy Mahjoory et al. (2019) and the neural syntax hypothesis Buzsáki (2010).

If we take each time step as 1 millisecond in the brain, then the refractory period  $\tau_\delta$  is around tens of milliseconds and the absolute refractory period  $\tau_r$  is around ten milliseconds. The coupling delay is around 50 millisecond Singer (2021). The integrative time window matches that of the coincidence detector (several millisecond König et al. (1996)). The frequency of oscillatory activity is around ten of milliseconds, within the Gamma band Tallon-Baudry & Bertrand (1999).

#### A.4.3 ABLATION STUDY OF THE TIMESCALE PARAMETERS

In this section, we provide ablation studies of time scale parameters in Table.2 (Fig.15).

Firstly, the coupling delay  $\tau_d = \tau_{d'}$  is most essential for the capability of the Composer. As shown in Fig.15, once removed, the parsing representation fails directly. As motivated by the neuroscientific studies on cortical dynamics Singer (2021), the delay coupling is also essential for generating the proper dynamical states in the Composer. How to understand this? As shown in Fig.2c,h, we need a chain of transient attractors to form a sequence oscillatory cell assemblies. Therefore, the positive feedback (the top-down feedback from DAE or from higher-level SCS) is desirable to wait for the modulated SCS neurons to recover from the refractory period. So that the cell assemblies sequence becomes a stable attractive states of the column dynamics. This idea is proved in Zheng et al. (2022) for a single column. In general, the feedback delay provide a “time window” for different metastable states to compete to emerge and disappear. It is notable that the effect of delay is soft and the Composer is robust to the shift of delay within a reasonable range (See Sensitivity Test).

Secondly, the refractory period ( $\tau_{r1}, \tau_{r2}$ ) has a secondary effect on the Composer. As shown in Fig.2c,h, the refractory period contributes to destroying the attractor dynamics (Fig.2b) to be transient



Figure 15: Ablation study of time scale parameters on all datasets.

(Fig.2c). If the refractory period is removed, the nested oscillatory states are not likely to be stable states any more. Therefore, all scores degrade for different degree.

Thirdly, the integration timescales ( $\tau_1, \tau_2, \tau_D$ ) also matters significantly. But the effect is relatively soft. As argued in Buzsáki (2010), the time window of readout neurons determines what they can “see”, which in turn determines what downstream modules (DAE) predict as top-down feedback. If the time window is very narrow, it is less plausible for the readout to “see” larger spatial-temporal scale coherence (or assemblies), which encodes higher-level objects. As a result, it is less likely to provide top-down predictive feedback for high-level objects. Instead, low-level objects are more likely to be captured by readout and top-down feedback. Taken together, the integration time window shapes the spatial-temporal scale of the cell assembly sequences and the time scale of the network dynamics (frequency of the oscillation).

Fourthly, the removal of the relative refractory period slightly degrades the nestedness of the Composer. The explanation is that: Representing the part-whole hierarchy is a combinatorial problem in nature, which needs to be iteratively “searched”. For example, when the object number increases as in the Ts dataset, the number of assemblies increases. Given larger number of assemblies as tree nodes, the possible configurations of the parse tree (correct or incorrect) gets exponentially larger. Therefore, the searching becomes harder because the searching space enlarges exponentially. However, while hard refractory period forces the system to switch among different states (spike fires at wrong timings), the ‘hardness’ could prevent efficient self-correcting once the system gets into a wrong state (because the hard refractory period constraints the available next firing timing). Thus, introducing a relative refractory period can help the system jump out of the local minimum, once it ‘finds’ much better states. It is likely that for this reason, enforcing  $g = 0$  in Fig.15 slightly degrades the Nest Score.

#### A.4.4 PARAMETER SENSITIVITY TEST

We provide a sensitivity test of parameters on SHOPS dataset in Fig.16. The score is relatively robust with the perturbation of the parameters, as long as the parameter is within the suitable range described by eq.28. Therefore, although including so many time scale hyper-parameters may be undesirable for machine learning models, these hyper-parameters only “softly” influence the performance of the Composer. Actually, we do not perform precise parameter tuning and these time scale parameters are “minimally” required for a brain-inspired model that takes time-scale interactions into account. In other words, we somehow need a time window to readout the spike in SCS, and that is all the Composer requires. Since the time scale parameters have a relatively wide effective range, the Composer is potentially generalizable to broadly bridge neuroscience and machine learning to study human-like vision in the future.



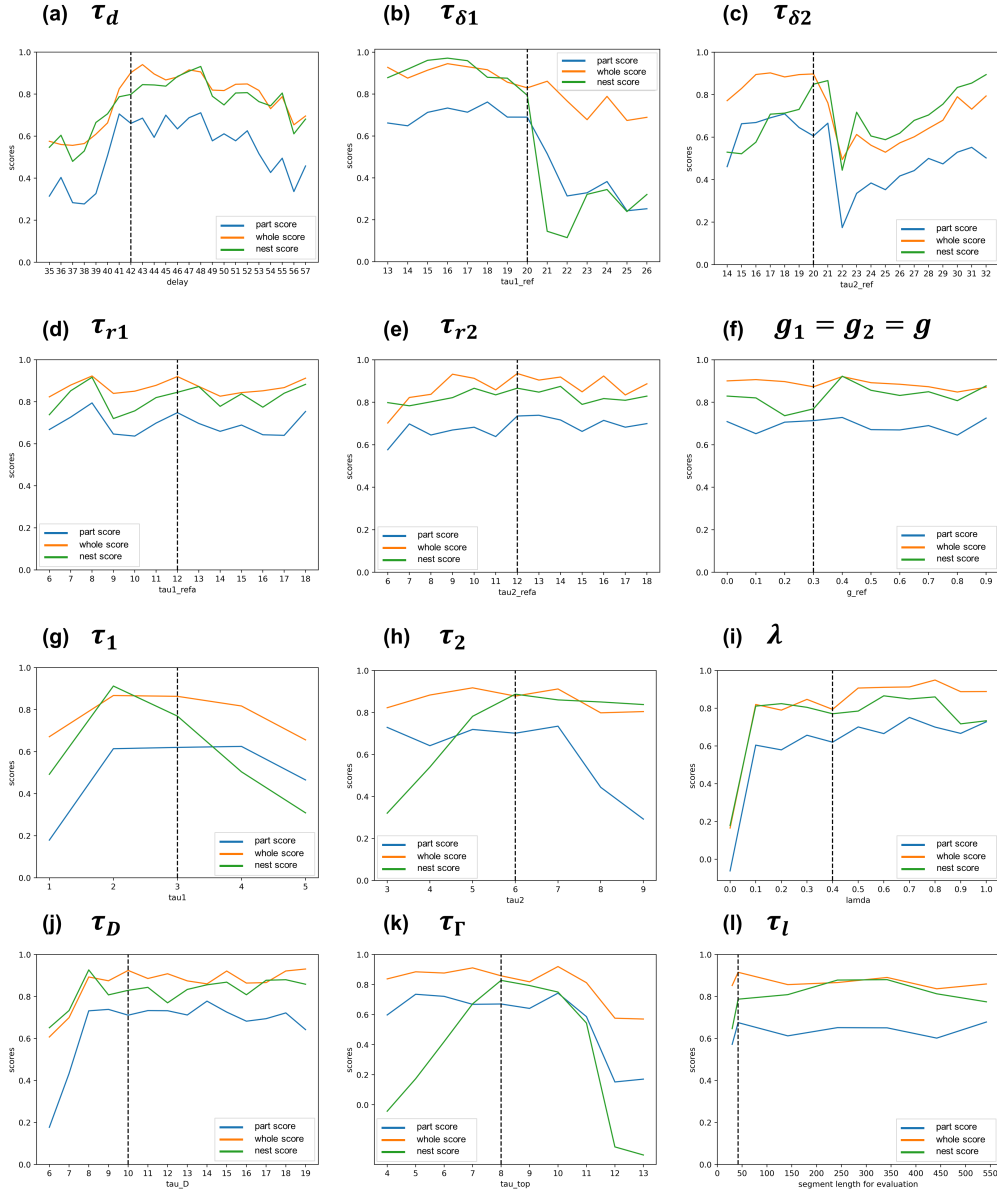


Figure 16: The sensitivity test of time scale parameters on SHOPS dataset. Black-dashed line indicates the value used for the Composer in the main text (Table.2). The value of parameters are perturbed to show how parsing degrades w.r.t parameter change. (a) The delay parameter of DAE feedback, same for part / whole level; (b) entire refractory period for part-level; (c) entire refractory period for whole level; (d) absolute refractory for part level; (e) absolute refractory for whole level; (f) the inhibitory effect of the relative refractory function; (g) the integration time window for part-level spiking neurons; (h) the integration window for whole-level spiking neurons; (i) the factor of the partial influence from skip connection; (j) the integration time window from part-level to whole level; (k) the integration time window from whole-level to part level; (l) the length of (spike train) segment used for evaluating the parsing quality (Table.1).

## A.5 TRAINING DETAILS

### A.5.1 RESOURCES

Our experiments have been performed on ubuntu 16.04.12 with devices: CPU (Intel(R) Xeon(R) CPU E5-2640 v4 @ 2.4GHz) and 4×GeForce RTX 2080 Ti. The python version is 3.6.3.

### A.5.2 NETWORK ARCHITECTURE AND TRAINING HYPERPARAMETERS

The details of training neural networks are shown in Table.4. All networks are trained with stochastic gradient descent (SGD).

Table 4: Details of training DAE

Dataset	encoder	decoder	learning rate	noise	minibatch size	epoch num
Ts (part)	FC(1600, 1000) Sigmoid()	FC(1000, 1600) Sigmoid()	1e-3	0.5	16	200
Ts (whole)	FC(1600, 1000) Sigmoid()	FC(1000, 1600) Sigmoid()	1e-3	0.5	16	200
Squares (part)	FC(3600, 400) Sigmoid()	FC(400, 3600) Sigmoid()	1e-3	0.8	16	200
Squares (whole)	FC(3600, 400) Sigmoid()	FC(400, 3600) Sigmoid()	1e-3	0.6	16	200
SHOPs (part)	FC(3600, 400) Sigmoid()	FC(400, 3600) Sigmoid()	1e-3	0.7	16	200
SHOPs (whole)	FC(3600, 400) Sigmoid()	FC(400, 3600) Sigmoid()	1e-3	0.7	16	200
Double-MNIST (part)	FC(6400, 2000) Sigmoid()	FC(2000, 6400) Sigmoid()	1e-3	0.5	16	200
Double-MNIST (whole)	FC(6400, 2000) Sigmoid()	FC(2000, 6400) Sigmoid()	1e-3	0.5	16	200

### A.5.3 DATASET FOR TRAINING DAE

The details of training dataset are shown in Table.5. Examples of the training data are visualized in Fig.17. The setting of DAE training dataset follows the convention in Zheng et al. (2022).

Table 5: Training dataset details

Dataset	Training size	Input dimension	Object number
Ts (part)	60000	40 × 40	1
Ts (whole)	60000	40 × 40	1
Squares (part)	60000	60 × 60	1
Squares (whole)	60000	60 × 60	1
SHOPs (part)	20000	60 × 60	1
SHOPs (whole)	20000	60 × 60	1
Double-MNIST (part)	60000	80 × 80	1
Double-MNIST (whole)	60000	80 × 80	1

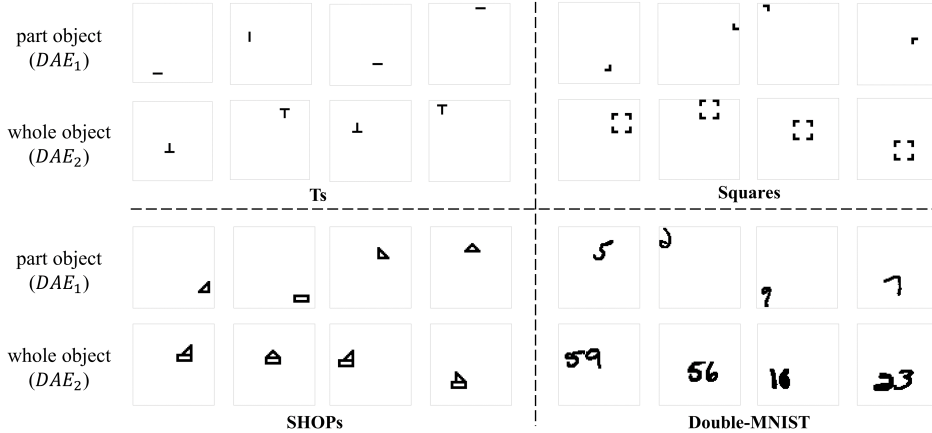


Figure 17: Examples of the training data to train part/whole level DAE.

#### A.5.4 LOSS FUNCTION

The DAE (either part or whole) are trained to minimize the MSE loss between the output of DAE and original image:

$$loss(x) = (x - DAE_i(\tilde{x}))^2, \quad i = 1, 2 \quad (29)$$

where  $x$  is the original single-object image in Fig.17.  $\tilde{x}$  is the corrupted version of  $x$ . Notably, the training of DAE has an unsupervised form and does not provide any explicit information on how to bind distributed features into the tree nodes or to form the parsing tree (coordination of multi-level tree nodes). These compositional structure all emerged during the simulation dynamics. All the training does is to provide the minimal prior about what (on average) the object (part/whole) looks like, so that the model could make sense of the multi-object scene at all. It is plausible that such priors of “object prototype” also exist in the brain to help parse the scene. For example, before parsing the house (Fig.1a in the main text), a person should have a prior about the door, window, and roof. Such prior should also influence the outcome of the parsing process.

In this paper, we treated these senses of the object (part or whole) as prior knowledge and explored how the parsing structure emerges on the condition of the prior knowledge. While some of the prior may be hard-wired in the brain through evolution, others may also be learned during development. The learning aspect of these priors is not discussed in this preliminary model, but recent work Zheng et al. (2023) provides insight into how it could potentially be achieved: the general architecture in this work, including the explicit separation of columns levels, the delay-coupled-bottom-up-top-down column architecture and hierarchical organization of time-scale constants, could be treated as the “**inductive bias**” of the end-to-end training. On the one hand, these biological constraints guarantee the nested states to be stable states of the network dynamics, given that the DAE is selective for different level objects (part/whole). On the other hand, the limited time window may also bias the DAE to “learn” to predict objects consistent with the scales of the time-window, because the time window softly determines what DAE can see “clearly” Zheng et al. (2023). Taken together, instead of training DAE separately, we could treat the entire model as a recurrent spiking neural network and train the model by back-propagation through time (BPTT) Wu et al. (2018). The loss function is only need to be modified minimally: the MSE loss between the entire scene (containing multiple whole-level objects) and the **averaged** top-down feedback (eg.  $\gamma_1(t)$ , or  $\Gamma(t)$ ). Due to the constraint of hierarchical temporal structure of the model (inductive bias), it is more efficient to learn a part-whole hierarchy representation to predict the whole image. Then, the single-object prior is possible to be learned in a fully unsupervised manner. We leave it as a promising future work: how neural syntax can be learned from predicting the external world, consistent with neuroscientific hypothesis Buzsáki et al. (2014).

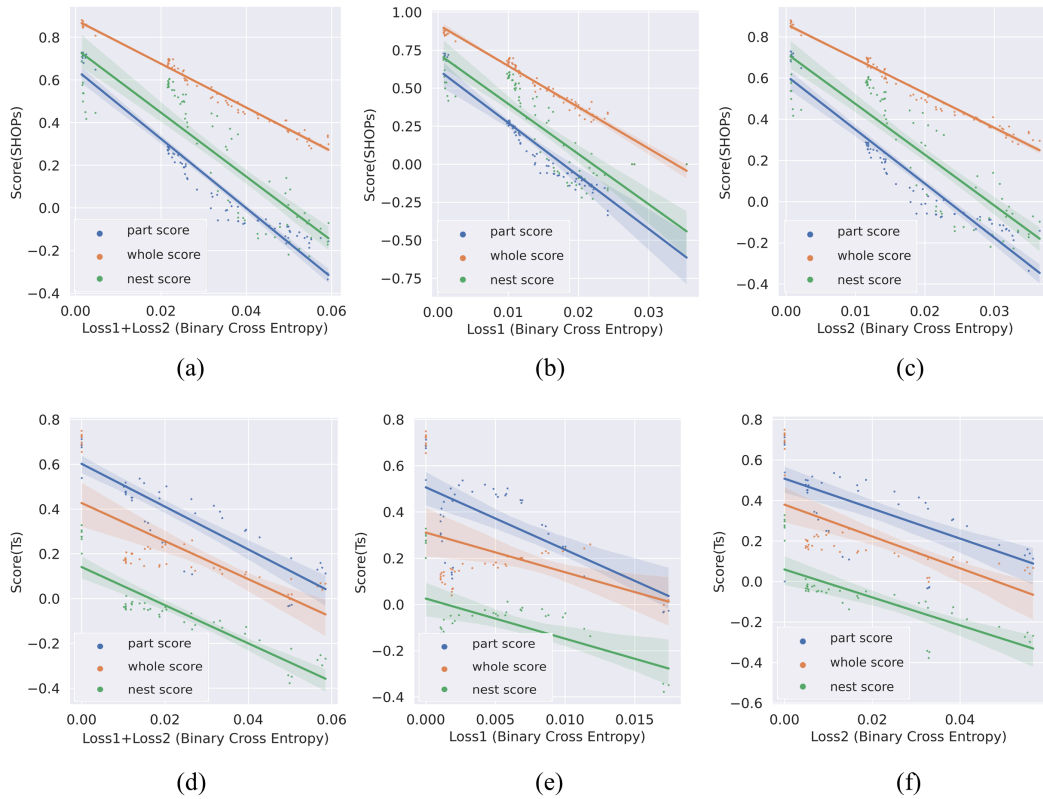


Figure 18: Loss vs Score. (a)(b)(c) results on SHOPs dataset; (d)(e)(f) results on Ts dataset. (a)(d) Relations between scores and total loss of part / whole level DAE ( $loss_1 + loss_2$ ); (b)(e) Relations between scores and total loss of part-level DAE ( $loss_1$ ); (c)(f) Relations between scores and loss of whole-level DAE ( $loss_2$ ); All results are consistent.

#### A.5.5 ABLATION STUDY OF THE DAE

Since denoised feedback from DAEs is an essential mechanism in the Composer, it is instructive to examine the relationship between the denoising performance of DAE and the scores. For this purpose, we randomly selected 100 learning rates from  $(10^{-3}, 1)$  and for each selected learning rate we trained one part-level DAE and one whole-level DAE. Then we evaluate the parsing ability of the Composer equipped with different DAEs. Fig.18 shows the relationship between the denoising loss and coherence scores. It is observed that lower loss positively correlates with higher scores on all metrics, indicating that there is direct interplays between denoising feedback attention and nested neuronal coherence. It make sense since the quality of DAE affects the quality of associative memory in Fig.2b and in turn affects the quality of coherence structure in Fig.2c,h,i.

## A.6 BIO-PLAUSIBILITY

In this section, we list and provide detailed discussion about the biological correlates of the design of the Composer in hope of inspiring future innovations.

1. **Delay-coupled oscillatory neural network:** In Singer (2021), the author describes the cerebral cortex as a delay-coupled recurrent oscillator network, which is very different from the architecture in the deep learning field. In the Composer, such architecture is captured, integrated within the deep learning framework, and acting as an essential ingredient of the mechanism (Section.A.4.3). The coupling delay provide a time window for alternative cell assemblies compete to emerge and disappear in order . In other words, the coupling delay makes the system non-Markovian and of infinite dimension (approximately, Izhikevich (2006)), so that the coding scheme and the capability of associative memory is much enlarged<sup>10</sup>
2. **Feed-forward and feed-back pathway along the cortical hierarchy:** In general, the cortex is organized into similar columns Douglas & Martin (2004), which is composed of six layers from layer I to layer VI (Fig.2d,g). Cortex are spatially organized corresponding to the spatial structure of the external physical world and hierarchically organized into levels. These basic features are captured in our model and act as essential elements for representation hypothesis: the representation of the part-whole hierarchy depends on such spatial and hierarchical organization. Notably, in Markov et al. (2013), the author also specifies the organization of the feedforward and feedback hierarchy. In detail, there are recognizable feedforward and feedback pathways between layer II/III of higher and lower level columns. This corresponds to the cross-level interaction between the part-level SCS and whole-level SCS in our model. Lastly, Markov et al. (2013) also shows that a long-distance feedforward path from lower level to high level exists in layer IIIb. These are realized as the skip connections from external driving input ( $x$ ) to the whole level SCS layer in the Composer (eq.4). More generally, the feedback from higher levels contains signals originating from both layer II/III and layer V/VI (corresponding to the latent space in the Composer). It is left to future work to study the “cross-level” interaction between pixel-level SCS and the latent space of DAEs.
3. **Time scale hierarchy:** Along the cortical hierarchy, there is a gradient of timescale hierarchy Mahjoory et al. (2019)—‘*We found that the dominant peak frequency in a brain area decreases significantly, gradually and robustly along the posterior-anterior axis, following the global cortical hierarchy from early sensory to higher order areas*’. Such time scale hierarchy is exploited in our model as the basis for representing hierarchical inclusion relationships among part-level cell assemblies and whole-level cell assemblies. However, since the frequency spectra are not unlimited, the capability of the part-whole representation may be limited by the range of the total frequency bands. Here, we treat such limitation as a shared weakness of our model and the brain, since the temporal resolution of cross-frequency coupling has shown to be a constraint for the capability of working memory of humans ( $7 \pm 2$ ) Nicola & Clopath (2019). Indeed, human also has a limited range of the hierarchy depth to represent instantaneously ( $\sim 5$  levels) Hinton (2021). At least three frequency bands could be explored in the future: gamma band, alpha band, and theta band.
4. **Topographical mapping:** As mentioned above, the spatial organization of the cortical column has a topographical correspondence to the physical world, called the topographical mapping Eickhoff et al. (2017). Such location-wise representation is exploited in GLOM Hinton (2021) as a core basis to represent the part-whole hierarchy and is similarly essential for our model. The topographical relationship enables the representation of objects as a set of (grouped) spatial regions with correlated activation patterns (grouped pixels with similar spike trains in the Composer or grouped columns with identical vectors, so called “identical islands of vectors” in GLOM). Besides, the location-wise representation also helps to clarify the inclusion relationship between whole and part across the hierarchy, both in our model and in GLOM. For example, the spatial region for part objects should also be spatially covered by their whole objects. However, spatial organization itself is not sufficient

<sup>10</sup>Each memory is realized as a trajectory instead of a single fixed point. The trajectory is the combination of transient fixed points so that the attractive states expand combinatorially or exponentially.

to represent part-whole hierarchy, because multiple co-activated features within the same level can lead to ambiguity (the binding problem Malsburg (1994)). As a result, we need similarity measure or coherence measure to specify “which is which”. This insight motivates the representation hypothesis of both spatial hierarchy in terms of topographical mapping and temporal hierarchy in terms of neuronal coherence.

5. **Abstract away the cortical BU/TD processing as autoencoder:** Predictive coding Rao & Ballard (1999) was first proposed by Dana H. Ballard and Rajesh P. N. Rao to explain the extra-classical receptive-field effects in primary visual cortex. Then, the predictive coding theory was mapped to the canonical circuit of cortical circuit Bastos et al. (2012) and served as a unified theory of brain function Friston (2010). In the predictive coding model, the bottom-up and top-down feedback attention (BU/TD) is formalized as the autoencoder architecture, and the reconstruction error should be minimized to achieve minimal “prediction error” or “surprise”. Such architecture is exploited in our model to realize the inner-column bottom-up / top-down pathways (BU/TD) and reconstruction error is minimized as the objective function of training. Interestingly, such predictive feedback is also related to the temporal coherence both in the cortex Engel et al. (2001) and in the Composer. Interestingly, “abstracting away certain part of a dynamical system in neuroscience as a learnable neural network” has the advantage of building-up more flexible dynamical models by training on large datasets. Therefore, it is plausible to generate new hypothesis beyond the traditional ones Eckstein et al. (2023); Peterson et al. (2021).
6. **Sparse code and dense code:** The dual coding scheme in the cortical circuits has been recognized when representing features: ultra-sparse coding in the superficial layer (layer II/III) and dense coding in deeper layer (layer V/VI) Tang et al. (2018); Wang (2018). While the latter encodes the statistical aspects of features, the former might additionally encodes the relationships. In this work, the dual coding scheme is realized as the sparse spike code in SCS and real-valued dense vector code in the DAE’s latent space. The synchrony in the SCS additionally encodes the relationship among objects.
7. **Relative refractory period:** Strictly speaking, the absolute refractory period (ARP) refers to the phase immediately after a spike initialization ( $\sim 2$  ms). The later phase where a spike is harder to be triggered (though not impossible) is referred to as the relative refractory period (RRP) Dayan & Abbott (2001). If we take 1-time step as 1 millisecond in real-world time, then the absolute refractory period is around 10 milliseconds (Table.2) in the Composer, which is much longer than the strict absolute refractory period. Therefore, the picture should be clarified as follows: the excitability of spiking neurons after a spike increases gradually, in the form of  $1 - e^{-t/\tau}$ . At the beginning phase, the excitability is low enough to prevent the neuron from firing a second spike given the conventional stimulus strength, but since the excitability increases rapidly during this phase, the relative period length is small compared to the whole refractory period. This beginning phase where the excitability is low enough compared to the stimulus strength but of fast increasing rate is treated as the absolute refractory period. In contrast, during the rest of the period, the excitability has recovered to the extent where neurons might generate a second spike but with a much lower probability. Since the recovery is much slower during the second phase, the temporal range is much longer than that of the first phase. This slow recovery phase is modeled as the relative refractory period in this work. The total refractory period can expand from tens of milliseconds to much longer, depending on the channel type on the axon of the neuron Gerstner et al. (2014). On the other hand, it is also conventional in numerical modeling that the absolute refractory period is modeled no less than 5 ms. In sum, the time scale of refractoriness fits the conventional setting of biological systems.
8. **Dendritic computation of pyramidal cell:** The driving signal and modulatory signal are distinguished in the cortical circuit Lee & Sherman (2010), where the driving signal acts on the proximal site of dendrites (near to the soma) and the modulatory signal acts on the distal sites (far from soma) Spruston (2008). The two types of inputs interact in a non-linear way. Such non-linear interaction between driving input and modulatory input is captured as the multiplication between the bottom-up integration and top-down modulation, realized as the pyramidal cells in the SCS (Fig.3ab in the main text). Such a gating effect inside the column is essential for the emergence of cell assemblies and the gating effect across levels is essential for the coordination of cell assemblies into nested temporal structure, so called nested neuronal coherence.

9. **Coincidence detector:** Abeles (1982) argued that cortical neurons in superficial layers are coincidence detectors, which detect sparse synchronous events within a narrow time window. In our model, the time constant of the integrative time window is small ( $\tau_1 \sim 2\text{ms}$ ,  $\tau_2 \sim 5\text{ms}$ ). As a result, the inner-level bottom-up integration of spiking activity in the superficial layer (pixel-wise SCS) is modeled as coincidence detectors. Such a narrow time window enables two things: (1) stochastic spikes fired at extremely adjacent time steps should be detected as a single event; (2) the temporal resolution of the synchronous event is kept within a small time-scale ( $\sim \tau_1, \tau_2$ ). Both are important to form a high-quality parse tree. Interestingly, a similar concept has also been developed in GLOM Hinton (2021), named as ‘coincidence filtering’.
10. **Meta-stability of cortical network:** “...*Single-trial analyses of ensemble activity in alert animals demonstrate that cortical circuit dynamics evolve through temporal sequences of metastable states. Metastability has been studied for its potential role in sensory coding, memory, and decision-making. Yet, very little is known about the network mechanisms responsible for its genesis...*” Mazzucato et al. (2015). In this work, we build such a system of metastable states (Fig.2) by integrating the spiking neural network (SNN) and artificial neural network (DAE as ANN) and further demonstrates its computational role in vision.
11. **Neuronal assembly as code words:** “*A widely discussed hypothesis in neuroscience is that transiently active ensembles of neurons, known as “cell assemblies,” underlie numerous operations of the brain, from encoding memories to reasoning. However, the mechanisms responsible for the formation and disbanding of cell assemblies and the temporal evolution of cell assembly sequences are not well understood...I suggest that the hierarchical organization of cell assemblies may be regarded as a neural syntax...*”Buzsáki (2010). Besides, assemblies are shown to be able to realize arbitrary computation function Papadimitriou et al. (2019). By combining machine learning models and neuroscientific constraints, we show how cell assembly can be transiently formed and disbanded, and be organized into a sequence at each level, and be hierarchically organized into spatial-temporal nested structure to express the neural syntax. More generally, various features, even of a continuous nature are represented as neuronal assemblies in the brain (population binary code), this coding scheme provides the basis to enable the Composer to deal with continuous features (RGB color) in the future Stockman (2019). The reservoir of neuronal assemblies could be more efficiently realized in neuromorphic devices Pei et al. (2019) to account for larger range of features. From the viewpoint of the Miehl et al. (2022), our model generate the assemblies by DAE-induced symmetry-breaking, which is one of the mechanism for assembly-generation.
12. **Temporal binding theory and feature integration theory:** Temporal binding theory Engel & Singer (2001); Malsburg (1994) and feature integration theory Wolfe (2020) are two mainstream theories to solve the binding problem: how distributed information is bound together to form the whole. The former is based on time coding and neuronal synchrony while the latter is based on top-down attention searching on a spatial map. The temporal synchrony, temporal coding, top-down attention, and spatial map are all captured in this model. Thus it is promising to explore whether it could serve as a canonical model to unify the two theories.
13. **The role of sequential / spatial attention:** The binding of distributed features or the segmentation of objects are also argued to be related to spatially-shift attentions Roelfsema (2023). In the Composer, such spatial / sequential attention co-emerges during the simulation (Fig.6e). Beyond theories that only highlight the role of attention, we further treat sequential attention and neuronal coherence as a “symbiotic couple”, and can be unified in a bigger picture of associative memory: both attention and neuronal coherence are partial description of certain-type associative memory (Fig.2). This bigger picture is consistent with recent findings Ramsauer et al. (2020), which suggests that associative memory and attention is actually the same thing, but viewed from different angles.
14. **Temporal-spatial theory of consciousness:** “*We postulate four different neuronal mechanisms accounting for the different dimensions of consciousness: (i) “temporospatial nestedness” of the spontaneous activity accounts for the level/state of consciousness as the neural predisposition of consciousness (NPC); (ii) “temporospatial alignment” of the pre-stimulus activity accounts for the content/form of consciousness as the neural prerequisite of consciousness (preNCC); (iii) “temporo-spatial expansion” of early stimulus-induced activity*



*accounts for phenomenal consciousness as neural correlates of consciousness (NCC); (iv) “temporo-spatial globalization” of late stimulus-induced activity accounts for the cognitive features of consciousness as the neural consequence of consciousness (NCCcon).” Northoff & Huang (2017). In this work, the temporospatial nestedness emerges and indicates the perceptual awareness (eg. recognizing the part-whole relationship), and the temporospatial alignment clarifies the content/form of the scene. The underlining assumption is that: the perceptual awareness emerges only if a part-whole relationship (of a visual scene or of a cognitive concept or of a plan or of a solution of certain problems) is well-organized in the internal representation of the brain (neural syntax). The quality of the representation reflect the “level of awareness” (eg. clear sight or dreamy sight). This assumption might be verified or falsified in future works.*

15. **Gamma oscillation and perceptual awareness:** *“...One theory suggests that rhythmic synchronization of neural discharges in the gamma band (around 40 Hz) may provide the necessary spatial and temporal links that bind together the processing in different brain areas to build a coherent percept. In this article we propose that this mechanism could also be used more generally for the construction of object representations that are driven by sensory input or internal, top-down processes...” Tallon-Baudry & Bertrand (1999). In this work, the spiking activity in the SCS approximately oscillates at the gamma band (tens of milliseconds if each time-step is regarded as 1 millisecond.) The gamma oscillation dynamically groups neurons into object representations (the representation hypothesis in the main text).*
16. **Preconfigured brain:** *In a recent “inside-out” conceptual framework of the brain, as Gyorgy Buzaki put it—“...This is the organization I call the preformed or preconfigured brain: a preexisting dictionary of nonsense words combined with internally generated syntactical rules. The neuronal syntax with its hierarchically organized rhythms determines the lengths of neuronal messages and shapes their combinations. Thus, brain syntax preexists prior to meaningful content...“Preconfigured” usually means experience-independent. The backbone of brain connectivity and its emerging dynamics are genetically defined. In a broader sense, the term “preconfigured” or “preexisting” is also often used to refer to a brain with an existing knowledge base, ....In the preconfigured brain model, learning is a matching process, in which preexisting neuronal patterns, initially nonsensical to the organism, acquire meaning with the help of experience...” Buzsáki (2019). Thus, the well-trained DAE in this paper could be treated as an essential preconfigured structure due to genetic codes or the life-long calibration of the sensory-action loop, which captures a range of object prototypes. Plasticity may only provide a secondary role to increase the precision of the ‘good-enough’ model Buzsáki (2019). Overall, the Composer is highly motivated by the Buzaki’s insights.*
17. **Plasticity:** *One of the designs that may depart from biology is that the connection weights are trained based on a gradient-based method instead of a correlation-based method, like Hebbian rule or spike timing plasticity Gerstner et al. (2014). However, this could be explained from two points of view. First, as argued above, the well-trained DAE could be regarded as the preconfigured structure which is gradually searched from evolution (amount to stochastic gradient-based search). Second, since the DAE structure in this model is relatively simple, the training objective (minimizing reconstruction error, the difference between input and feedback) could be interpreted as increasing the correlation between sensory neurons and modulatory neurons, so that the gradient-based training equals correlation-based plasticity. Indeed, Melchior & Wiskott (2019) shows that gradient-based learning and Hebbian plasticity can be unified in case of a simple autoencoder. Further, we could imagine that there is a two-stage learning algorithm, like the wake-sleep cycle: during the day, the system infers entities based on learned weight, during the night, the learned objects replay and the system efficiently updates the weight by association, which corresponds to the training phase of the DAE. Similar treatment has also been discussed in GLOM Hinton (2021).*
18. **Inner-layer recurrent connection:** *Another design feature that may depart from the biological brain is that the spike coding space (SCS) itself is not recurrent in our model. However, this could also be explained from at least two points of view. First, the feedforward and recurrent connection usually have different functional roles in the cortical circuit, and have different levels of domination. For example, layer IV in the visual cortex are*

mainly feedforward and the recurrent effect are relatively weak. As a result, the inner-layer recurrence of SCS are treated as secondary compared to the recurrence of inner-column top-down feedback or inter-level top-down feedback. So that it is temporally ignored for simplicity. Further, the localized inner-level recurrence may play a secondary role (different from that of top-down feedback) to speed up the convergence by forming a grid frame (by spatially-organized connection) to encode the prior of the proximity property of objects (Gestalt principles Wagemans et al. (2012)). Secondly, the entire two-layer column could be recognized as a single layer, with DAE parameterizing the recurrent connection weight among spiking neurons, similar to Dmitry Krotov’s idea in Kozachkov et al. (2023). And the general mechanism still works. In other words, there is no restriction to view the two-level system as a column or a layer. In either case, the models maintain their bio-plausibility.

19. **Polychronization** refers to the generalization of absolute synchronization into structured asynchrony. As argued in Izhikevich (2006), due to the heterogeneity and conduction delay of the neural system, polychrony is more plausible than absolute synchrony. While the externally observed spike firing time is asynchronous, the arriving time of asynchronous spikes to downstream readout neurons is (internally) synchronous. In other words, the shift in spiking time is balanced out by the shift in conduction delay. According to Buzsáki (2010), the more rigorous definition of cell assemblies should be based on internal observation (downstream readout neurons) instead of external observation (human observer). Therefore, polychronous representation is in essence also synchronous representation. In Composer, due to the cross-level coupling-delay the part-level cell assemblies and whole-level cell assemblies seems to have a fixed temporal shift (Table.1). However, from a readout neuron perspective (eg. SCS spiking neurons), the spike arrival of whole level feedback ( $\gamma$ ) is just in time to enslave the spike firing of part-level SCS neurons. For this reason, we ignored the slight temporal shift when computing the score (Section.A.3.6). Besides, to make the visualization more intuitive, we compensate the temporal shift for visualization in Fig.6, Fig.7 in the main text: synchronization is more intuitive to capture the nestedness than polychronization, although they all indicates the coherence state.

## A.7 DETAILS ON EXPERIMENTS AND ADDITIONAL RESULTS.

### A.7.1 CONVERGENCE

In Fig.8, we show the convergence of scores along the iteration. 100 randomly selected samples are used, and the score are evaluated every 100 time steps (so  $3000/100=30$  data point in total for each score).

Convergence on the SHOPs dataset achieves the best overall results. However, the potential overlap of part-level objects when composing the whole-level object imposes additional challenges on recognizing the part-level objects, indicated by the relatively lower Part Score in Fig.8a.

Convergence on the Squares dataset is very interesting. On the one hand, the Whole Score takes the lead all the time, indicating that global information is firstly recognized by the Composer, which is very similar to human vision Lee & Nguyen (2001) and is also consistent with Gestalt psychology. On the other hand, the Part Score and Nest Score undergo an initial descending period before going up. Here, we explain this phenomenon: Compared with the very starting phase, where spikes are randomly and densely fired, the emergence of whole-level squares around  $500 \sim 1000$  time steps provide new conditions on the part level. While this has benefits in the long run, it could degrade the representation in the short run, because the Composer needs to rethink its representation and make modifications. For example, the part-level firing becomes sparser, and there are more incorrect synchronizations. This may degrade the part-level grouping and hierarchical coordination (nestedness). In other words, the fact that each whole object is composed of four parts complicates the self-correcting / searching process, after the whole-level objects are recognized. Fortunately, after a short period of self-correcting, the Scores go up again and gradually converge to expected neuronal coherence as in other cases.

Convergence on the Ts dataset is very challenging due to the object number is much larger. On the one hand, the Composer needs to distinguish 6 wholes and 12 parts. On the other hand, 6 whole objects and 12 part objects impose  $6^{12}$  potential combinations of part-whole relationships (each part can choose to belong to one of the six wholes). Therefore, it takes time to search for / sample the optimal configuration. Even if the parts/wholes are grouped by neuronal coherence, there is a high possibility that the parts and wholes are not well coordinated to form proper nested structure. Since the neural computation in the brain can also be regarded as sampling Buesing et al. (2011), these challenges may also cause problems in perception like the binding problem Engel & Singer (2001); Von der Malsburg (1999). Therefore, on the Ts dataset, the Nest Score lags behind the other scores.

Convergence on Double-Digit MNIST is also challenging for the Composer because the objects are of much higher diversity. Therefore, it is harder for the Composer to clearly distinguish the objects and to form well-synchronized cell assemblies. Therefore, the Part Score is lower than other scores and the variance is higher than in other cases. However, it is surprising that the Composer still achieves good parsing, indicated by the convergent Nest Score, even though objects are less recognizable.

### A.7.2 VISUALIZATION

In Fig.6 in the main text, we visualize the spiking pattern, attention map, and local field potential along the convergent process. To better visualize the cell assemblies, we reorder the index of neurons on the y-axis (Fig.6c) so that neurons encoding the same object are close on the y-axis. Besides, in order to distinguish different cell assemblies more clearly, we color the spikes based on the ground truth assignment of the neurons in SCS (= pixels in the image), so that the color of the cell assembly indicates what object the cell assembly represents. In other words, the represented objects can be directly recognized by comparing the color of the cell assembly and ground truth. This fact can be verified by comparing the circled cell assemblies in Fig.6c, the circled zoomed-in spike patterns in Fig.6d, and the circled objects in the ground truth (Fig.6a). It is clear that the synchronized cell assemblies gradually emerges from randomness along the simulation. Each synchronized cell assembly represents the parts/wholes of the object. Cell assemblies at different levels are coordinated properly according to the part-whole relationship.

Surprisingly, the emergent parsing tree reform itself during the dynamics in Fig.6 c: From phase II to phase III, the part assemblies (colored by deep and light green) reversed their order! This observation demonstrates the multi-stable property of the part-whole solution in the Composer, which is a dynamical system in nature. This property is essential to account for the uncertainty and diversity

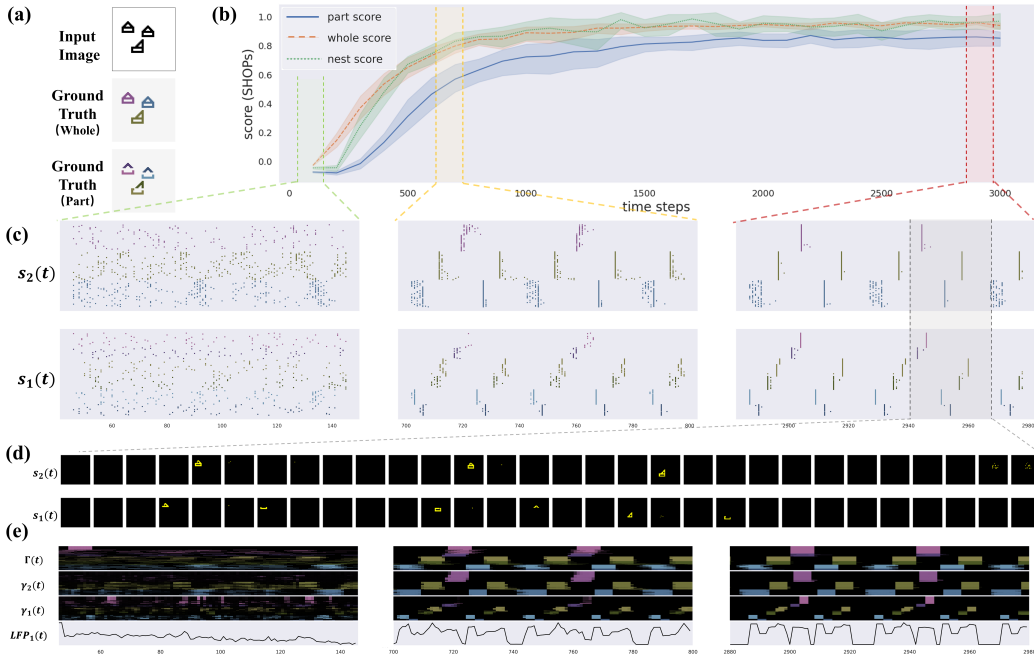


Figure 19: Additional visualization on SHOPs dataset.

of the part-whole relationship and challenges most traditional artificial neural networks with only feedforward connection or only by supervised training Greff et al. (2020).

In Fig.6e, it is also observed that different types of top-down attention also emerge into structured patterns. To keep consistent with Fig.6c, the neuron indexes are also reordered and the attention map is also colored based on the ground truth. The depth of the color reflects the value of the attention map. It is observed that the structured pattern has the same order as the spiking pattern, yet of long timescales. This indicates that attention plays a role in modulating the spike timing in SCS. However, such modulation is not single-way, but a iterative interplay between bottom-up integration and top-down modulation. Therefore, both DAE and SCS play essential roles in solving the parsing problem.

In Fig.6, we also shows the emergence of the oscillatory LFP at the part level, which is the summed top-down feedback:  $LFP_1(t) = \sum_{i=1}^N \gamma_{1i}(t)$ , where  $i$  is the neuron index in the part level.

### A.7.3 MORE VISUALIZATIONS

In Fig.7, we briefly show the visualization results on other datasets. Here we provide more detailed visualization results on the four datasets Fig.19 to Fig.26. Two cases are provided for each dataset, including one normal case and one fail case.

We also provide a zip file containing videos to visualize temporal evolution of neuronal activities in SI, about 60MB.

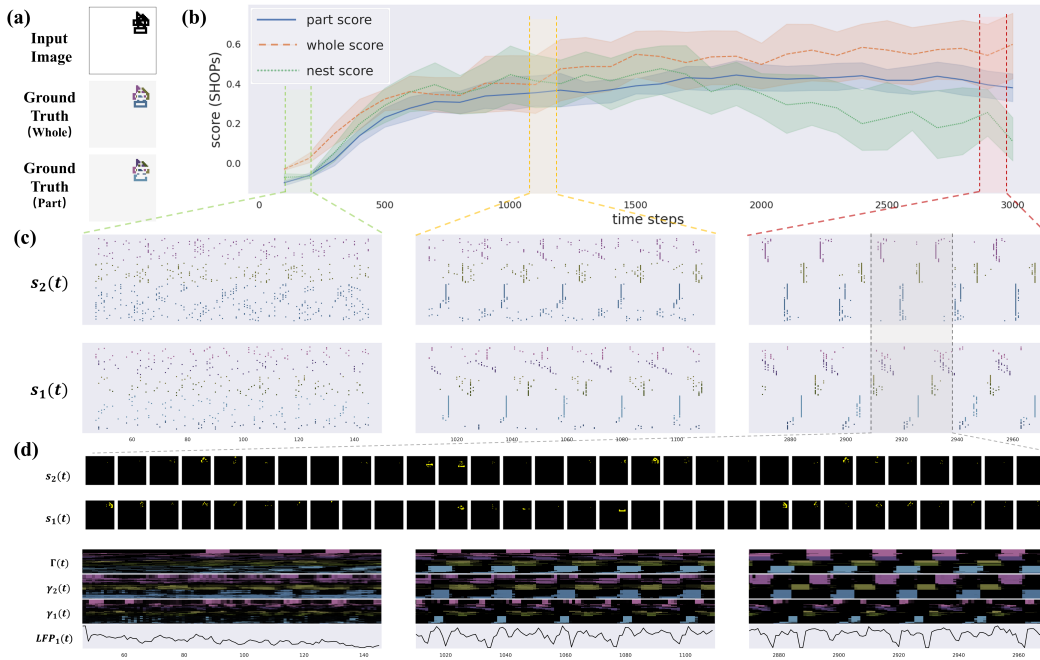


Figure 20: Additional visualization on SHOPS dataset: The fail case, when objects sickly overlap.

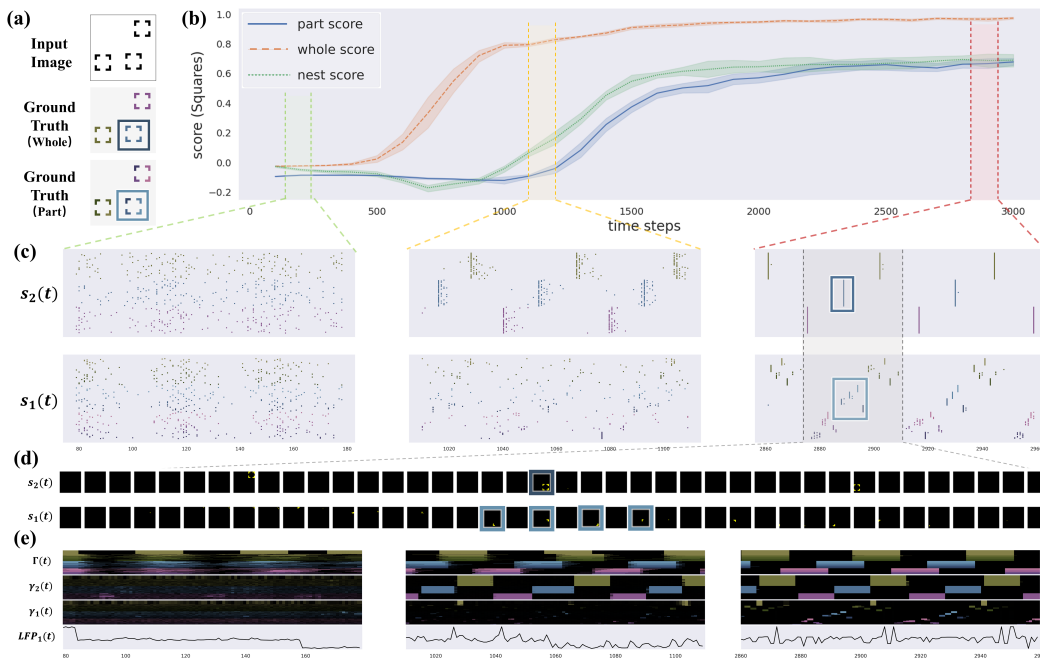


Figure 21: Additional visualization on Squares dataset: Squares are not overlap.

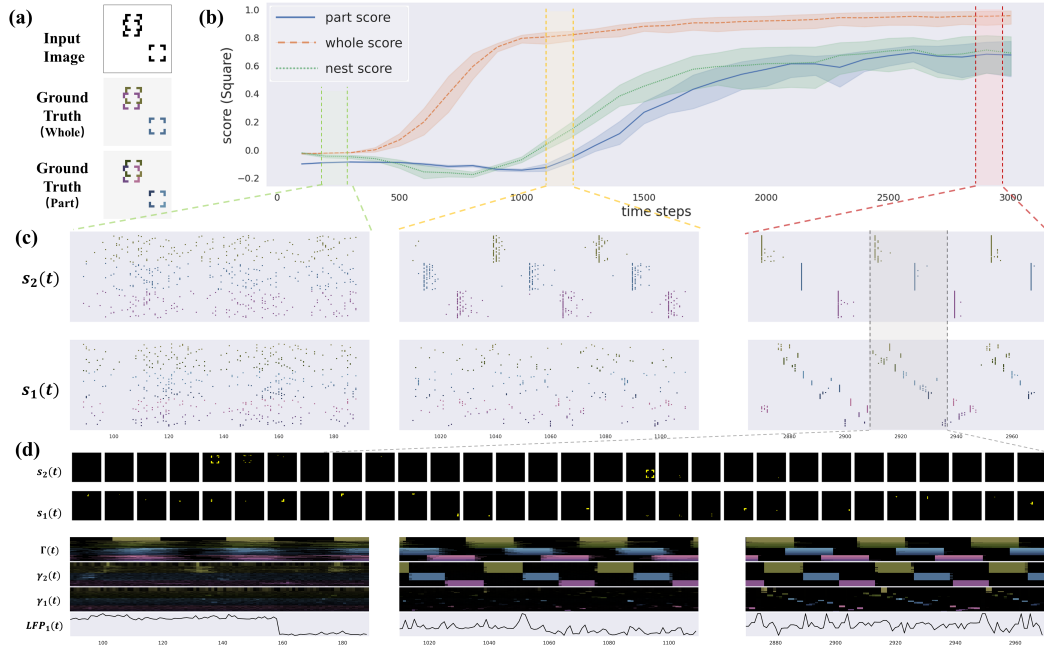


Figure 22: Additional visualization on Squares dataset: Two Squares heavily overlap but the nested structure remains.

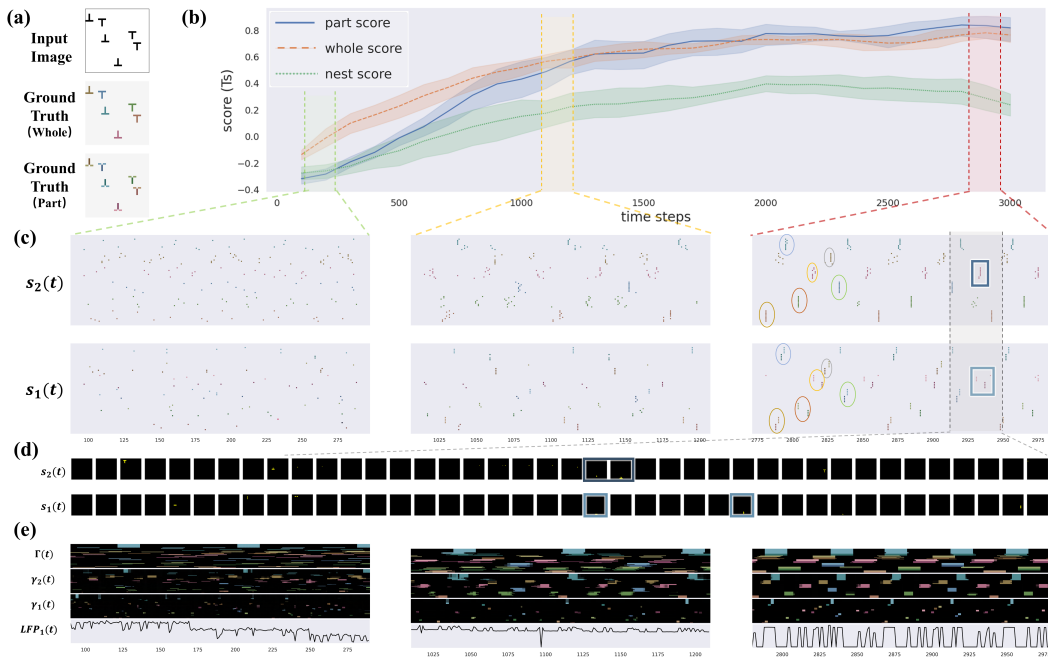


Figure 23: Additional visualization on Ts dataset. Colored circled indicates the coordinated cell assemblies. Same color indicates the part-whole relationship.

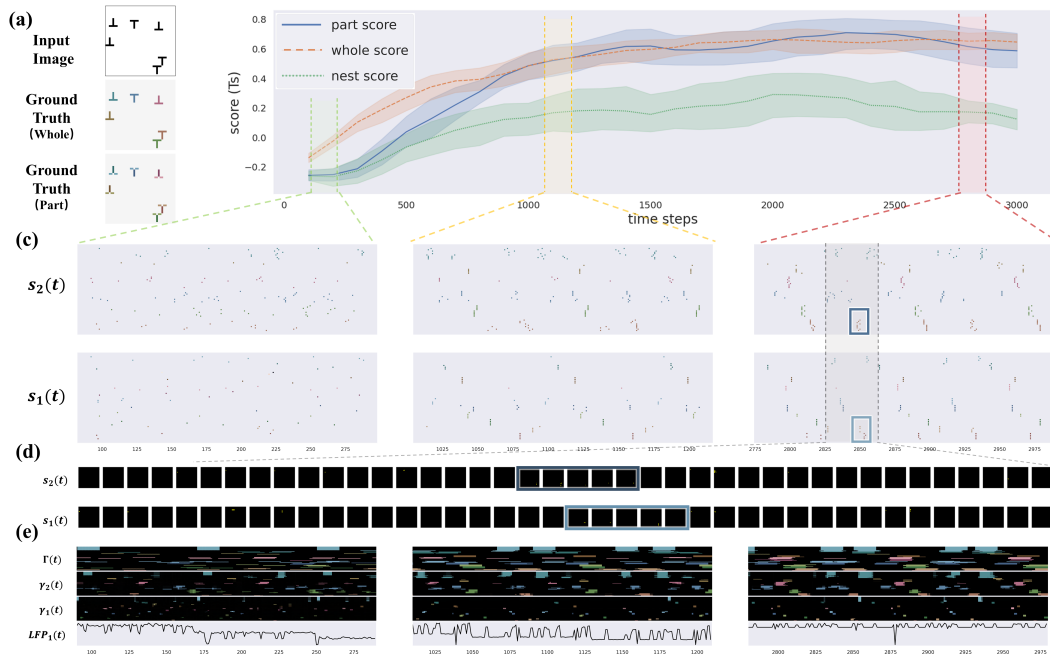


Figure 24: Additional visualization on Ts dataset. Nestedness is not as clear as fig.23.

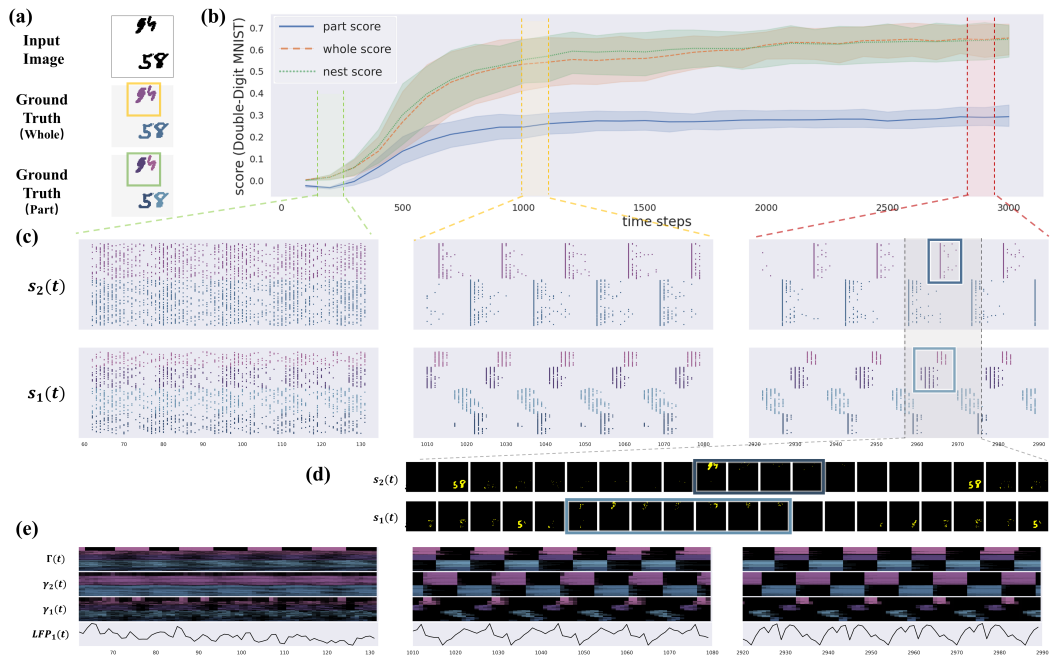


Figure 25: Additional visualization on Double-Digit-MNIST dataset.

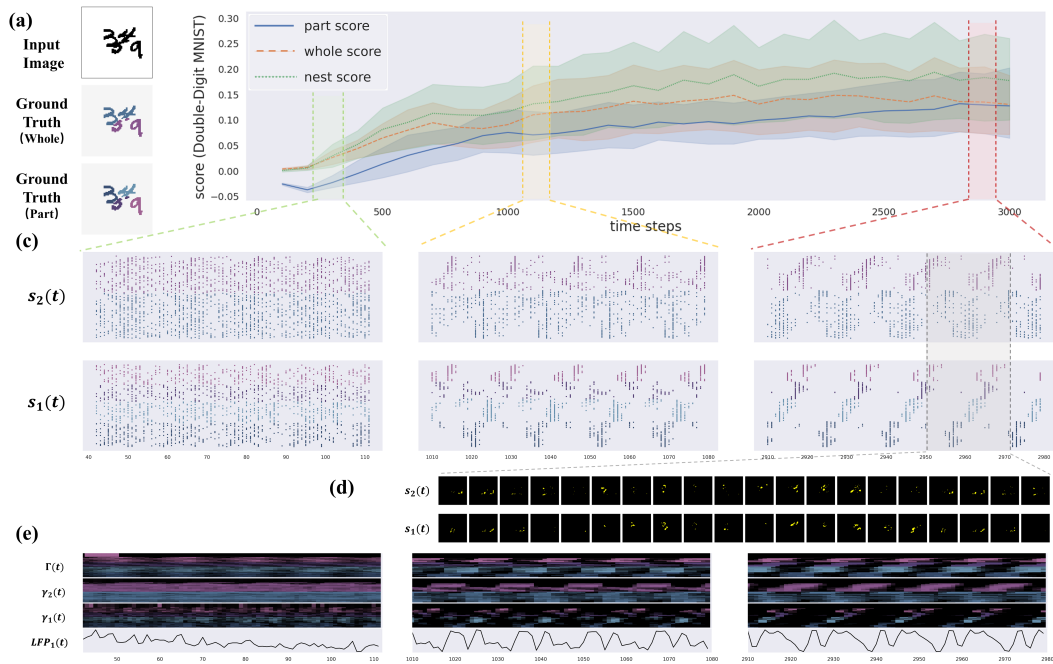


Figure 26: Additional visualization on Double-Digit-MNIST dataset: Digits are crowded and the coherence structure is weaker.