

THE FUNCTION OF CONTEXTUAL ILLUSIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Many visual illusions are contextual by nature. In the orientation-tilt illusion, the perceived orientation of a central grating is repulsed from or attracted towards the orientation of a surrounding grating. An open question in vision science is whether such illusions reflect basic limitations of the visual system, or if they correspond to corner cases of neural computations that are efficient in everyday settings. Here we develop deep recurrent network architectures that approximate neural circuits linked to contextual illusions (Mély et al., 2018). We show that these architectures, which we refer to as γ -Nets, are more sample efficient for learning contour detection than the state of the art, and exhibit an orientation-tilt illusion consistent with human data. Correcting this illusion significantly reduces γ -Net performance by driving it to prefer low-level edges over high-level object boundary contours. Overall, our study suggests that contextual illusions are a byproduct of neural circuits that help biological visual systems achieve robust and efficient perception, and that incorporating such circuits in artificial neural networks can improve computer vision.

1 INTRODUCTION

An open debate since the inception of vision science concerns why we experience visual illusions. Consider the class of “contextual” illusions, for which the perceived qualities of an image region, such as its orientation or color, are systematically altered by its surrounding. A well-studied contextual illusion is the orientation-tilt illusion depicted in Fig. 1a, where the perception of a central grating is influenced by a grating of a different orientation in the surround (O’Toole & Wenderoth, 1977). When the two orientations are similar, the central grating appears tilted slightly away from the surround (repulsion; Fig. 1a, top). When the two orientations are dissimilar, the central grating appears tilted slightly towards the surround (attraction; Fig. 1a, bottom). Does the contextual bias of the orientation-tilt illusion reflect a bug or a feature of neural computations?

Over the past 50 years, there has been a number of neural circuit mechanisms proposed to explain individual contextual illusions (reviewed in Mély et al., 2018). Recently, Mély et al. (2018) theorized a unified cortical circuit, constrained by physiology of primate visual cortex (V1), that explains contextual illusions across visual domains – from the orientation-tilt illusion to color induction. Illusions arise in this circuit from recurrent interactions between neurons that tile retinotopic visual space, leading to complex contextual (center/surround) effects. For the orientation-tilt illusion, separate neural populations encoding the surround grating’s orientation can suppress or facilitate neural population encodings of the central grating’s orientation. Asymmetric competition between these two modes of interaction causes surround repulsion to predominate when the center and surround are similar, and surround attraction to predominate when they are dissimilar.

The neural circuit of Mély et al. (2018) explains how contextual illusions might emerge, but it does not explain why. One possibility is that contextual illusions like the orientation-tilt illusion are “bugs”: vestiges of evolution or a by-product of biological constraints on the neural hardware. Another possibility is that contextual illusions are “features”: the corollary of efficient neural routines for scene segmentation (Keemink & van Rossum, 2016; Mély et al., 2018). Here, we provide computational evidence for the latter possibility and demonstrate that the orientation-tilt illusion reflects neural strategies optimized for segmentation.

Contributions We introduce the γ -Net, a trainable and hierarchical implementation of the neural circuit of Mély et al. (2018). We demonstrate that the γ -Net (i) is more sample efficient than state-

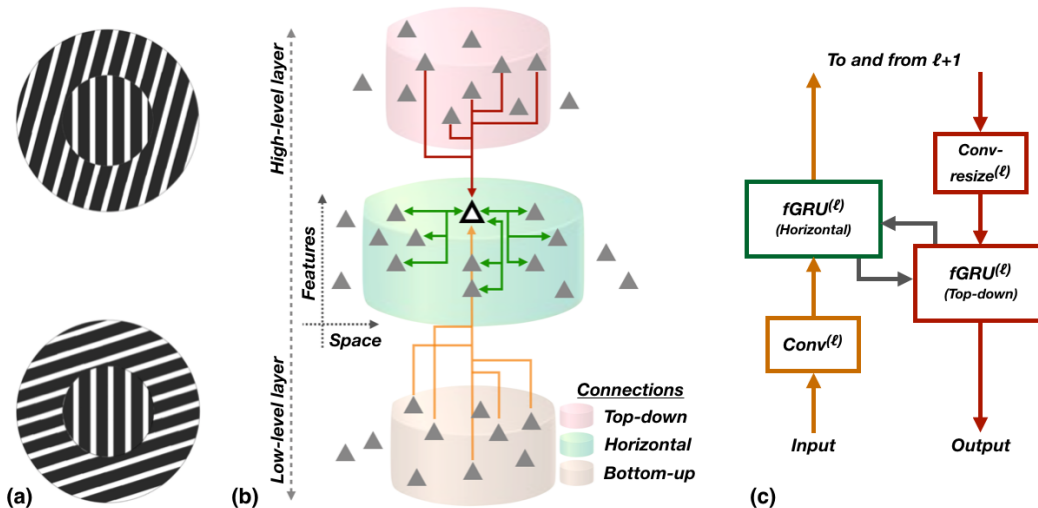


Figure 1: **The orientation tilt-illusion (O’Toole & Wenderoth, 1977) is a contextual illusion where a central grating’s perceived orientation is influenced by a surround grating’s orientation.** (a) When a central grating has a similar orientation to its surround, it is judged as tilting away (repulsed) from the surround. When the two gratings have dissimilar orientations, the central grating is judged as tilting towards (attracted) the surround. (b) We extend the recurrent circuit proposed by Mély et al. (2018) to explain this and other contextual illusions into a hierarchical model that can simulate horizontal (within a layer) and top-down (between layers) interactions between units. Triangles depict units at different spatial positions and encoding different features. Arrows describe the direction of interactions between units. A unit of interest (outlined triangle) is influenced by instantaneous bottom-up (orange), recurrent top-down (red) and horizontal (green) connections. (c) A deep network schematic of the diagram in (b), which forms the basis of the γ -Net introduced here. Horizontal and top-down connections are implemented with feedback gated recurrent units (fGRUs). Image encodings sweep through these blocks on every timestep, from bottom-up (left path) to top-down (right path), and predictions are read out from the fGRU closest to image resolution on the final timestep. This motif can be stacked (e.g., extending to $\ell + 1$).

of-the-art convolutional architectures on two separate contour detection tasks, and (ii) exhibits an orientation-tilt illusion after it is optimized for contour detection, which we show emerges from its preference for high-level object-boundary contours over low-level edges.

2 RELATED WORK

Modeling the visual system Convolutional neural networks (CNNs) are often considered the *de facto* “standard model” of vision. CNNs and their extensions represent the state of the art for most computer vision applications with performance approaching – and sometimes exceeding – human observers on certain visual recognition tasks (He et al., 2016; Lee et al., 2017; Phillips et al., 2018). CNNs also provide the best fit to rapid neural responses in the visual cortex (see Kriegeskorte 2015; Yamins & DiCarlo 2016 for reviews). Nevertheless, multiple lines of evidence suggest that biological vision is still far more robust and versatile than CNNs (see Serre, 2019, for a recent review). CNNs suffer from occlusions and clutter (Fyall et al., 2017; Rosenfeld et al., 2018; Tang et al., 2018). They are also sample inefficient at learning visual relations (Kim et al., 2018) and solving simple grouping tasks (Linsley et al., 2018b). State-of-the-art CNNs require massive datasets to reach their impressive accuracy (Lake et al., 2015) and their ability to generalize beyond training data also appears to be somewhat more limited than originally anticipated (Geirhos et al., 2018; Recht et al., 2018).

It is well established that cortical feedback contributes to the robustness of biological vision (Hochstein & Ahissar, 2002; Wyatte et al., 2014; Kafaligonul et al., 2015). Feedforward projections in the visual system are almost always matched by feedback projections (Felleman & Van Essen, 1991), and feedback has been implicated in visual “routines” that cannot be implemented through purely

feedforward vision, such as incremental grouping or filling-in (O’Reilly et al., 2013; Roelfsema, 2006). There is also a growing body of work that is demonstrating the potential of recurrent neural networks (RNNs) to account for neuroscience data (Fyall et al., 2017; Klink et al., 2017; Siegel et al., 2015; Tang et al., 2018; Nayebi et al., 2018; Kar et al., 2019; Kietzmann et al., 2019).

Feedback for computer vision In contrast to CNNs, which build processing depth through a cascade of filtering and pooling stages, RNNs learn to dynamically process stimuli by re-using filtering operations over “timesteps” of recurrence. RNNs were originally developed to process temporal sequences (e.g., Mozer 1992), but their benefits have also been demonstrated even for the processing of static images. Notable efforts include multi-dimensional RNNs (Graves et al., 2007), which treat images as pixel sequences, and convolutional-RNNs, which learn spatial kernels to compute RNN activities. There are now many successful applications of RNNs in computer vision, including object recognition and super-resolution tasks (Liang & Hu, 2015; Kim et al., 2016).

In the current work, we are motivated by the horizontal gated recurrent unit (hGRU, Linsley et al. 2018a), which approximates the recurrent neural circuit model of (Mély et al., 2018) for explaining contextual illusions. The hGRU was designed to learn a difficult synthetic incremental grouping task, and a single layer of the hGRU learned long-range spatial dependencies that CNNs with orders-of-magnitude more weights could not. The hGRU relaxed some of the biological constraints from the original circuit model – including the assumption of non-negativity, which guaranteed separate and competing processing stages for suppression vs. facilitation. In our extension of the hGRU we re-introduce these constraints, which are in principle necessary for explaining contextual illusions, and extend it into a hierarchical formulation that can model horizontal connections (between units within a layer) and top-down connections (from units in a higher layer to units in a lower layer), to compete with state-of-the-art hierarchical models for contour detection in naturalistic images.

3 METHODS

A neural circuit model for contextual illusions We begin by introducing the dynamical neural circuit of Mély et al. (2018), which explains contextual illusions by simulating interactions between cortical hypercolumns tiling the visual field. In the model, hypercolumns are indexed by their 2D coordinate (x, y) and feature channels k . Units in hypercolumns encode idealized responses for a visual domain (e.g., idealized neural responses from the orientation domain were used to simulate the orientation-tilt illusion). Dynamics of a single unit at xyk obey the following equations (we bold activity tensors to distinguish them from learned kernels and parameters):

$$\begin{aligned} \eta \dot{H}_{xyk}^{(S)} + \epsilon^2 H_{xyk}^{(S)} &= \left[\xi Z_{xyk} - (\alpha H_{xyk}^{(F)} + \mu) C_{xyk}^{(S)} \right]_+ & \# \text{ Stage 1: Recurrent suppression of } \mathbf{Z} \\ \tau \dot{H}_{xyk}^{(F)} + \sigma^2 H_{xyk}^{(F)} &= \left[\nu C_{xyk}^{(F)} \right]_+, & \# \text{ Stage 2: Recurrent facilitation of } \mathbf{H}^{(S)} \end{aligned}$$

where

$$\begin{aligned} C_{xyk}^{(S)} &= (W^S * \mathbf{H}^{(F)})_{xyk} & \# \text{ Compute suppression interactions} \\ C_{xyk}^{(F)} &= (W^F * \mathbf{H}^{(S)})_{xyk}. & \# \text{ Compute facilitation interactions} \end{aligned}$$

Circuit activities consist of a feedforward drive, recurrent suppression, and recurrent facilitation, respectively denoted as $\mathbf{Z}, \mathbf{H}^{(S)}, \mathbf{H}^{(F)} \in \mathbb{R}^{X \times Y \times K}$ (X is width, Y is height of the tensor, and K is its feature channels). The circuit takes its “feedforward encodings” \mathbf{Z} from hypercolumns (e.g., orientation encodings from hypercolumn units), and introduces recurrent suppressive and facilitatory interactions between units, $\mathbf{C}^{(S)}, \mathbf{C}^{(F)} \in \mathbb{R}^{X \times Y \times K}$. These interactions are implemented with separate kernels for suppression and facilitation, $W^S, W^F \in \mathbb{R}^{E \times E \times K \times K}$, where E is the spatial extent of connections on a single timestep (constrained by primate physiology*). Because these interactions

*Mély et al. (2018) use different connectivity patterns for short vs. long-range suppression and facilitation. Following the derivation of Linsley et al. (2018a), we simplify notation by summarizing these connections as separate kernels for suppression vs. facilitation.

are implemented through convolutions, they serially spread over timesteps of processing to connect units positioned at different spatial locations. The circuit outputs $\mathbf{H}^{(F)}$ after reaching steady state.

The key assumption of the circuit for asymmetric suppression vs. facilitation is implemented with separate stages for computing suppression vs. facilitation. Hidden states $\mathbf{H}^{(S)}$ and $\mathbf{H}^{(F)}$ are updated with suppression vs. facilitation interactions, respectively. Suppression, but not facilitation, is applied directly to the recurrent output of the circuit. Thus, given a fixed amount of suppression and facilitation, suppression increases alongside unit activities, whereas excitation remains constant.

Circuit integration, suppression, and facilitation are controlled by hand-tuned parameters. Linear and multiplicative suppression (i.e., shunting inhibition) are controlled by scalars μ and α , feedforward drive is modulated by the scalar ξ , and linear facilitation is controlled by the scalar ν . Circuit time constants are scalars denoted by η , ϵ , τ and σ . All activities are non-negative and both stages are linearly rectified (ReLU) $[\cdot]_+ = \max(\cdot, 0)$.

Feedback gated recurrent units Linsley et al. (2018a) developed a version of this circuit for computer vision applications, called the hGRU, which they trained with gradient descent to fit its connectivity and parameters to image datasets (unlike the original circuit, for which these parameters were tuned by hand). The hGRU replaced time constants with dynamic gates, converted the suppressive recurrent state $\mathbf{H}^{(S)}$ into an instantaneous activity, and introduced a term for quadratic facilitation. These changes were made to improve performance on a synthetic contour task, and for their purposes, it was not important to maintain constraints for contextual illusions or building a hierarchical model (they used single-layer models).

We extend the hGRU formulation in two key ways. First, we enforce non-negativity. The circuit of Mély et al. (2018) explains contextual illusions like the orientation-tilt illusion because suppression is applied to feedforward encodings \mathbf{Z} before facilitation, and the strength of surround suppression but not surround facilitation is modulated by the magnitude of the circuit output. The non-negativity constraints we introduce into the fGRU are necessary to guarantee separate stages of suppression followed by facilitation that can implement asymmetric contextual interactions. Second, we extend the circuit into a feature processing hierarchy, and develop separate configurations for implementing horizontal connections between units within a layer of a network, and top-down connections between units in different layers of a network. We call our module the feedback gated recurrent unit (fGRU). The following equations describe the evolution of fGRU recurrent units in $\mathbf{H} \in \mathbb{R}^{X \times Y \times K}$, which are influenced by the non-negative feedforward encodings $\mathbf{Z} \in \mathbb{R}^{X \times Y \times K}$ (e.g., a convolutional layer’s response to a stimulus) over discrete timesteps, denoted by $\cdot[t]$:

Stage 1:

$$\begin{aligned} \mathbf{G}^S &= \text{sigmoid}(U^S * \mathbf{H}[t-1]) && \# \text{ Compute channel-wise selection} \\ \mathbf{C}^S &= W^S * (\mathbf{H}[t-1] \odot \mathbf{G}^S) && \# \text{ Compute suppression interactions} \\ \mathbf{S} &= \left[\mathbf{Z} - \left[(\alpha \mathbf{H}[t-1] + \mu) \mathbf{C}^S \right]_+ \right]_+, && \# \text{ Suppression of } \mathbf{Z} \end{aligned}$$

Stage 2:

$$\begin{aligned} \mathbf{G}^F &= \text{sigmoid}(U^F * \mathbf{S}) && \# \text{ Compute channel-wise recurrent updates} \\ \mathbf{C}^F &= W^F * \mathbf{S} && \# \text{ Compute facilitation interactions} \\ \tilde{\mathbf{H}} &= \left[\nu(\mathbf{C}^F + \mathbf{S}) + \omega(\mathbf{C}^F * \mathbf{S}) \right]_+ && \# \text{ Facilitation of } \mathbf{S} \\ \mathbf{H}[t] &= (1 - \mathbf{G}^F) \odot \mathbf{H}[t-1] + \mathbf{G}^F \odot \tilde{\mathbf{H}}. && \# \text{ Update recurrent state} \end{aligned}$$

Like the original circuit, the fGRU has separate stages for suppression (\mathbf{S}) and facilitation (\mathbf{H}). In the first stage, the feedforward encodings \mathbf{Z} are suppressed by non-negative interactions between units in $\mathbf{H}[t-1]$ (an fGRU hidden state from the previous timestep). Suppressive interactions are computed with the kernel $W^S \in \mathbb{R}^{E \times E \times K \times K}$, where E describes the spatial extent of horizontal connections on a single timestep (see Appendix A for details on how this is set in γ -Nets). This

kernel is convolved with a gated version of the persistent hidden state $\mathbf{H}[t - 1]$. The gate activity \mathbf{G}^S is computed by applying a sigmoid nonlinearity to a convolution of the kernel $U^S \in \mathbb{R}^{1 \times 1 \times K \times K}$ with $\mathbf{H}[t - 1]$, which transforms its activity into the range $[0, 1]$. Linear and multiplicative forms of suppression are controlled by the parameters $\mu, \alpha \in \mathbb{R}^K$, respectively.

In the second stage, linear and multiplicative facilitation is applied to the instantaneous activity \mathbf{S} . The kernels $W^F \in \mathbb{R}^{E \times E \times K \times K}$ controls facilitation interactions. Additive and multiplicative forms of facilitation are scaled by the parameters $\nu, \omega \in \mathbb{R}^K$, respectively. A gate activity is also computed during this stage to update the persistent recurrent activity \mathbf{H} . The gate activity \mathbf{G}^F is computed by applying a sigmoid to a convolution of the kernel $U^F \in \mathbb{R}^{1 \times 1 \times K \times K}$ with \mathbf{S} . This gate updates $\mathbf{H}[t]$ by interpolating $\mathbf{H}[t - 1]$ with the candidate activity $\tilde{\mathbf{H}}$. After every timestep of processing, $\mathbf{H}[t]$ is taken as the fGRU output activity. As detailed in the following section, the fGRU output hidden state is either passed to the next convolutional layer (Fig. 1c, $\text{fGRU}^{(\ell)} \rightarrow \text{conv}^{(\ell+1)}$), or used to compute top-down connections (Fig. 1c, $\text{fGRU}^{(\ell+1)} \rightarrow \text{fGRU}^{(\ell)}$).

The fGRU can learn connections between units either within a layer or between layers (Fig. 1b). These two configurations stem from changing the activities used for a fGRU’s feedforward encodings and recurrent hidden state. “Horizontal connections” between units within a layer are learned by setting the feedforward encodings \mathbf{Z} to the activity of a preceding convolutional layer, and setting the hidden state \mathbf{H} to a persistent activity initialized as zeros (Fig. 1c, $\text{conv}^{(\ell)} \rightarrow \text{fGRU}^{(\ell)}$). “Top-down connections” between layers are learned by setting fGRU feedforward encodings \mathbf{Z} to the persistent hidden state $\mathbf{H}^{(\ell)}$ of a fGRU at layer ℓ in a hierarchical model, and the hidden state \mathbf{H} to the persistent activity $\mathbf{H}^{(\ell+1)}$ of an fGRU at a layer one level up in the model hierarchy (Fig. 1c, $\text{fGRU}^{(\ell+1)} \rightarrow \text{fGRU}^{(\ell)}$). The functional interpretation of the top-down fGRU is that it first suppresses activity in the lower layer using the higher layer’s recurrent horizontal activities, then, and then applies a kernel to the residue for facilitation, which allows for computations like interpolation, sharpening, or “filling in”. Note that an fGRU for top-down connections does not have a persistent state (it mixes high and low-level persistent states), but an fGRU for horizontal connections does.

γ -Net Our main objective is to test how a model with the capacity for contextual illusions performs on natural image analysis. We do this by incorporating fGRUs into leading feedforward architectures for contour detection tasks, augmenting their “bottom-up” processing with modules for learning horizontal and top-down connections (Fig. 1c). We refer to the resulting hierarchical models as γ -Nets, because information flows in a loop that resembles a γ : on every timestep, image encodings make a full bottom-up to top-down cycle through the architecture, and dense predictions are read-out from the lowest-level recurrent layer of the network (thus, information flows in at the top of the hierarchy, loops through the network, and flows out at the top of the hierarchy). In our experiments we convert leading architectures for two separate contour detection problems into γ -Nets: A VGG16 for BSDS500 (He et al., 2019), and a U-Net for detection of cell membranes in serial electron microscopy images (Lee et al., 2017).

4 CONTOUR DETECTION EXPERIMENTS

Overview We evaluated γ -Net performance on two contour detection tasks: object contour detection in natural images (BSDS500 dataset; Martin et al., 2004) and cell membrane detection in serial electron microscopy (SEM) images of mouse cortex (Kasthuri et al., 2015) and mouse retina (Ding et al., 2016). Different γ -Net configurations were used on each task, with each building on the leading architecture for their respective datasets. All γ -Nets use 8-timesteps of recurrence and instance normalization (the latter of which to minimize the vanishing gradient problem typical of RNNs; Ulyanov et al., 2016; Coijmans et al., 2017, see Appendix A for details). The γ -Nets were trained with Tensorflow and NVIDIA Titan RTX GPUs using single-image batches and the Adam optimizer (Kingma & Ba, 2014, dataset-specific learning rates are detailed below). Models were trained with early stopping. Training was stopped if the validation loss did not drop for 50 straight epochs, and the weights with the best validation-set performance were used for test.

Model evaluation We evaluated models in two ways. First, we validated them against state-of-the-art models for each contour dataset using standard benchmarks. As discussed below, we verified that our implementations of these state-of-the-art models matched published performance. Second,

we tested sample-efficiency after training on subsets of the contour datasets without augmentations. Measuring sample-efficiency compares the inductive biases of different architectures, and is critical for understanding how the capacity for exhibiting contextual illusions influences performance. Model “wall time” (i.e., how long training took) is outside the scope of the current work, as it is a function of hardware/software optimizations and does not necessarily reflect the relative advantages of different inductive biases. Model performance is evaluated as F1 ODS after post-processing (Martin et al., 2001; 2004), as is standard for contour detection tasks.

4.1 OBJECT CONTOUR DETECTION IN NATURAL IMAGES

Dataset We trained models for object contour detection on the BSDS500 dataset. The dataset contains object-contour annotations for 500 natural images, which are split into train (200), validation (100), and test (200) sets.

Architecture details The leading approach to BSDS500 is the Bi-Directional Cascade Network (BDCN, He et al. 2019), which places multi-layer readout modules at every processing block in a ILSVRC12-pretrained VGG16, and optimizes a loss that balances contributions from each of these readouts to achieve better scale tolerance in final prediction. All leading deep learning approaches to BSDS500 begin with a VGG16 pretrained on ILSVRC12 object recognition (He et al., 2019).

Our γ -Net for BSDS500 begins with the same ILSVRC12-pretrained VGG16. fGRUs were introduced for learning horizontal (conv2_2, conv3_3, conv4_3, conv5_3) and top-down connections (conv5_3 \rightarrow conv4_3, conv4_3 \rightarrow conv3_3, and conv3_3 \rightarrow conv2_2). To pass top-down activities between layers, higher-level activities were resized to match lower-level ones, then passed through two layers of 1×1 convolutions with linear rectification, which registered feature representations from higher-to-lower layers. The γ -Net was trained with learning rates of $3e^{-4}$ on its randomly initialized fGRU weights and $1e^{-5}$ on its VGG-initialized weights.

In contrast to the BDCN (and other recent approaches to BSDS) with multiple read-outs and engineered loss functions, we take γ -Net predictions as a linear transformation of the lowest fGRU layer in its feature hierarchy, and optimize the model with binary cross entropy between per-pixel predictions and labels (following the method of Xie & Tu, 2017)). This approach works because the γ -Net executes bottom-up to top-down processing sweeps on every timestep. The ability to segment object contours means that the model uses feedback to merge low- and high-level image feature representations at the bottom of its feature hierarchy, resembling classic “V1 scratchpad” hypotheses for computation in visual cortex (Gilbert & Sigman, 2007; Lee & Mumford, 2003). We compared γ -Nets with a BDCN implementation released by the authors, which was trained using the routine described in (He et al., 2019)[†].

Results We validated the γ -Net against the BDCN after training on a full and augmented BSDS training set (Xie & Tu, 2017). The γ -Net performed similarly in F1 ODS (0.802) as the BDCN (0.806) and humans (0.803), and outperformed all other approaches to BSDS (Fig. 2a; Deng et al. 2018; Xie & Tu 2017; Hallman & Fowlkes 2015; Kokkinos 2015; Wang et al. 2019; Liu et al. 2019).

Our hypothesis is that contextual illusions reflect routines for efficient scene analysis, and that the capacity for exhibiting such illusions improves model sample efficiency. Consistent with this hypothesis, the γ -Net was up to an order-of-magnitude more efficient than the BDCN. A γ -Net trained on 5% of BSDS performs on par with a BDCN trained on 10% of the BSDS, and a γ -Net trained on 10% of the BSDS performs on par with a BDCN trained on 100% of BSDS. Unlike the BDCN, the γ -Net trained on 100% of BSDS outperformed the state of the art for non-deep learning based models (Hallman & Fowlkes, 2015), and nearly matched the performance of the popular HED trained with augmentations (Xie & Tu, 2017). We also compared the γ -Net to versions of the models with lesions that tested the importance of its (i) horizontal connections, (ii) top-down connections, (iii) timesteps of processing, and (iv) weight sharing over timesteps of processing. The full γ -Net outperformed each of these lesioned versions on every subset of BSDS500 (Fig. S4).

[†]We replicated published results with this implementation. In our experiments, we train the BDCN with the same regularization, batch size, and optimizer as He et al. (2019). For sample efficiency experiments, we searched through multiple learning rates and found that it had no affect on performance on these small BSDS500 datasets (Fig. S3).

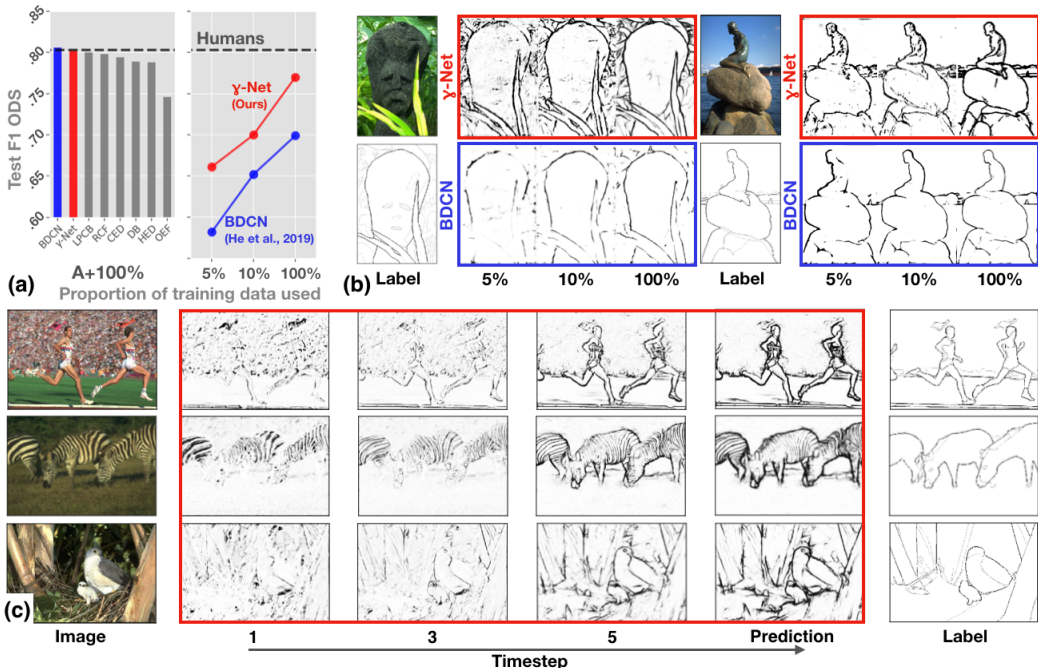


Figure 2: **Object contour detection in BSDS500 images.** (a) The γ -Net is on par with humans and the state-of-the-art for contour detection (BDCN; He et al. 2019) when trained on the entire training dataset with augmentations. In this regime, it also outperforms the published F1 ODS of all other approaches to BSDS500 (LPCB: Deng et al. 2018, RCF: Liu et al. 2019, CED: Wang et al. 2019, DB: Kokkinos 2015, HED: Xie & Tu 2017, and OEF: Hallman & Fowlkes 2015). The γ -Net outperforms the BDCN when trained on 5%, 10%, or 100% of the dataset. Performance is reported as F1 ODS (Martin et al., 2001). (b) BDCN and γ -Net predictions after training on the different proportions of BSDS500 images. (c) The evolution of γ -Net predictions across timesteps of processing. Predictions from a γ -Net trained on 100% of BSDS are depicted: its initially coarse detections are refined over processing timesteps to select figural object contours.

We examined recurrent feedback strategies learned by γ -Net for object contour by visualizing its performance on every timestep of processing. This was done by passing its activity at a timestep through the final linear readout layer. The γ -Net iteratively refines its initially coarse contour predictions. For example, the top row of Fig. 2c shows that the γ -Net selectively enhance the boundaries around the runner’s bodies while suppressing the feature activities created by the crowd. In the next row of predictions, salient zebra stripes are gradually suppressed in favor of body contours.

4.2 CELL MEMBRANE DETECTION

Datasets Deriving a connectome of the brain by mapping the connections between neurons is an important step towards understanding the algorithms that they implement (Briggman & Bock, 2012). CNNs can automate this task by detecting neuron membranes in every pixel/voxel of high-resolution serial electron microscope (SEM) images. Challenges like SNEMI3D (Kasthuri et al., 2015), which contains annotated SEM images of mouse cortex, have helped drive progress towards automation. Here, we test models on membrane detection in SNEMI3D and a separate SEM dataset of mouse retina (“Ding” from Ding et al. 2016). We split both datasets into training (80 images for SNEMI3D and 307 images for Ding) and test sets (20 images for SNEMI3D and 77 images for Ding). Next, we generated versions of each training dataset with 100%, 10%, or 5% of the images, as well as versions of the full datasets augmented with random left-right and up-down flips (A+100%).

Architecture details The current state-of-the-art on SNEMI3D is a variant of the popular U-Net architecture (Ronneberger et al., 2015), which uses a different depth, different number of feature maps

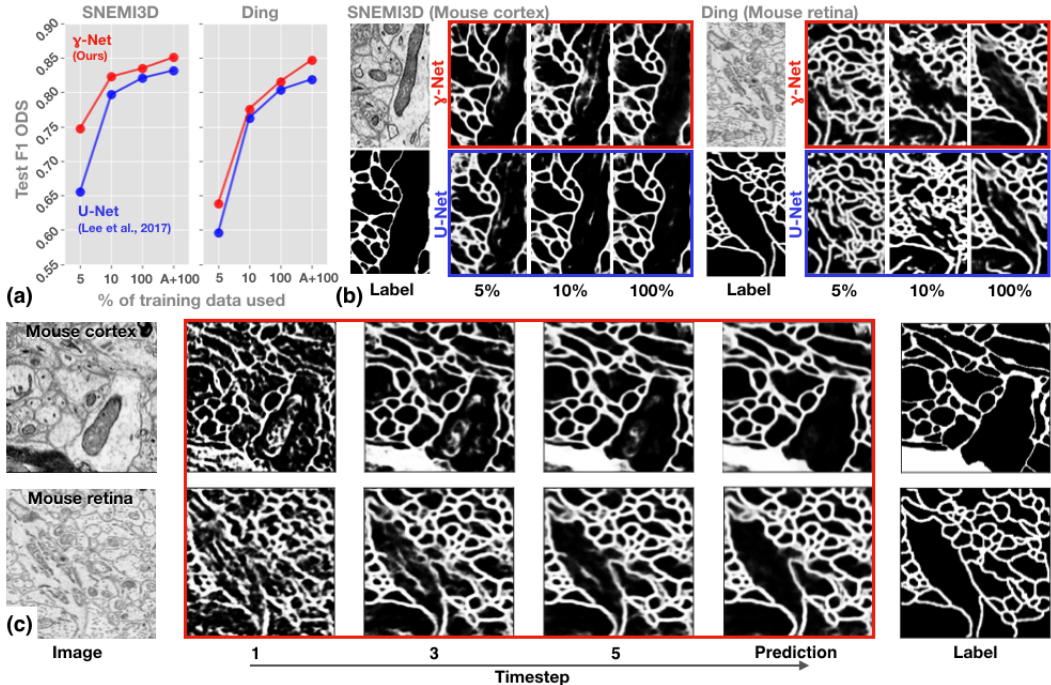


Figure 3: **Membrane prediction in serial electron microscopy (SEM) images of neural tissue.** (a) The γ -Net outperforms a state-of-the-art U-Net (Lee et al., 2017) for membrane detection when trained on SEM image datasets of mouse visual cortex and retina tissue. Performance is F1 ODS (Martin et al., 2001). (b) Network predictions after training on different proportions of each dataset. (c) The evolution of γ -Net predictions across timesteps of processing after training on 100% of the datasets. γ -Net learns to iteratively suppress contours belonging to internal cell features, such as organelles, which should not be annotated as contours for the purpose of neural tissue reconstruction.

at every layer, and introduces additional operations such as residual connections (Lee et al., 2017). We developed a γ -Net for cell membrane segmentation that contains recurrent connections implemented by fGRUs. The U-Net of Lee et al. (2017) contains four encoder blocks (convolution, pooling and subsampling) and four decoder blocks (transpose convolution and convolution). The γ -Net replaces the convolutional residual layers within each of these blocks with a single layer of convolutions followed by an fGRU (as in the high level diagram of Fig. 1c, $\text{Conv}^{(\ell)} \rightarrow \text{fGRU}^{(\ell)}$). In the encoder pathway, fGRUs store horizontal interactions between spatially neighboring units of the preceding convolutional layer in their hidden states. In the decoder pathway, fGRUs learn top-down connections that connect recurrent units from higher-feature processing layers to lower-feature processing ones, modifying encoder fGRU recurrent states.

Like in Lee et al. (2017), this γ -Net uses spatial pooling in its encoder pathway, and transpose convolution in its decoder pathway to upsample pooled activities and enable top-down interactions between layers of the hierarchy. The model’s final prediction is linearly read-out from the first fGRU layer, which maintains a feature representation of the same height and width as the input image. These γ -Nets were trained from a random initialization with a learning rate of $1e^{-2}$ to minimize class-balanced per-pixel binary cross-entropy. We verified our implementation of the U-Net from Lee et al. (2017) by demonstrating ”superhuman” segmentation performance in cell segmentation on SNEMI3D (Appendix B).

Results The γ -Net and U-Net of Lee et al. (2017) performed similarly when trained on full augmented versions of both the SNEMI3D and Ding datasets (Fig. 3a A+100%). However, γ -Nets were consistently more sample efficient than U-Nets on every reduced dataset condition (Fig. 3b).

We visualized the recurrent membrane detection strategies of γ -Nets trained on 100% of both datasets. Membrane predictions were obtained by passing neural activity at every timestep through the final

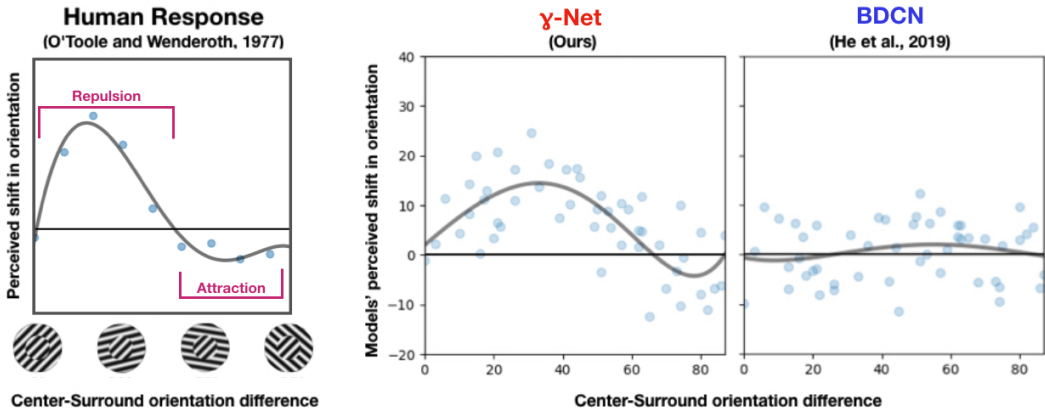


Figure 4: **Optimizing for contour detection produces an orientation-tilt illusion in the γ -Net .** The orientation-tilt illusion (O’Toole & Wenderoth, 1977) describes how perception of the center grating’s orientation is repulsed from the surround when the two are in similar orientations (e.g., $\approx 30^\circ$), and attracted to the surround when the two are in dissimilar (but not orthogonal) orientations (e.g., $\approx 60^\circ$). We test for the orientation-tilt illusion in models trained on BSDS500 contour detection. Model weights were fixed and new layers were trained to decode the orientation of grating stimuli of a single orientation. These models were tested on grating stimuli in which surround orientations were systematically varied w.r.t. the center (exemplars depicted in the left panel). The γ -Net but not the BDCN had an orientation-tilt illusion. Gray curves depict a fourth-order polynomial fit.

linear readout. The γ -Net prediction timecourse indicates that it learns a complex visual strategy for membrane detection: it gathers a coarse “gist” of membranes in the first timestep of processing, and iteratively refines these predictions by enhancing cell boundaries and clearing out spurious contours of elements like cell organelles (Fig. 3c).

5 ARE ILLUSIONS FEATURES OR BUGS?

Orientation-tilt illusion Like the neural circuit of Mély et al. (2018), fGRU modules allow units to asymmetrically suppress and facilitate each other. This potentially gives γ -Nets (which contain fGRU modules) the capacity to exhibit similar contextual illusions as humans. Here, we tested whether a γ -Net trained on contour detection in natural images exhibits an orientation-tilt illusion.

We trained orientation decoders on the outputs of models trained on the full BSDS500 dataset. These decoders were trained on 100K grating images, in which the center and surround orientations were the same (Fig. S2a). These images were sampled from all orientations and spatial frequencies. The decoders had two 1×1 convolution layers and an intervening linear rectification to map model output into the sine and cosine of grating orientation. Both the γ -Net and BDCN achieved nearly perfect performance on a held-out validation set of gratings.

We tested these models on a set of 1K grating stimuli generated with different center-surround grating orientations (following the method of O’Toole & Wenderoth 1977, Fig. S2b), and recorded model predictions for the center pixel in these images (detailed in Appendix C). Surprisingly, γ -Net encodings of these test images exhibited a similar tilt illusion as found in human perceptual data (Fig. 4b), with repulsion when the central and surround gratings had similar orientations, and attraction when these gratings were dissimilar. This illusory phenomenon cannot be explained by accidental factors such as aliasing between the center and the surround, which would predict the opposite pattern, indicating that the illusion emerges from the model’s strategy for contour detection. In contrast, the BDCN which only relies on feedforward processing to detect contours did not exhibit the effect (Fig. 4b). We also tested for the orientation-tilt illusion in lesioned γ -Nets, but found that both repulsion and attraction regimes are only present when the full model is intact (Fig. S4).

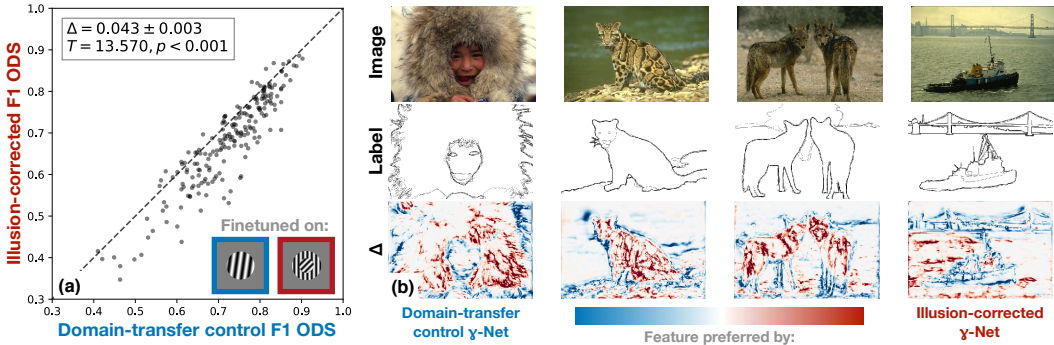


Figure 5: **Contour detection performance of the γ -Net depends on an orientation-tilt illusion.** (a) F1 ODS scores on BSDS500 test images (200 total) from γ -Nets after correcting an orientation-tilt illusion (“illusion-corrected”) or not (“domain-transfer control”). The domain-transfer control γ -Net was trained to decode the orientation of single-grating stimuli (blue), and the illusion-corrected γ -Net was trained to decode the orientation of the central grating in illusory grating stimuli (red). Readouts for decoding orientation were fixed, and γ -Net weights were allowed to change during training. Per-image F1 ODS was significantly greater for The domain-transfer control γ -Net than the illusion-correction γ -Net . (b) The illusion-corrected γ -Net was biased towards low-level contours, whereas the domain-transfer control γ -Net was biased towards contours on object boundaries.

Correcting the orientation-tilt illusion What strategies does the orientation-tilt illusion reflect? We tested this question by measuring performance on BSDS500 while we corrected the orientation-tilt illusion of a γ -Net . A γ -Net that displayed the orientation-tilt illusion was trained to decode the central grating orientation of tilt-illusion stimuli (Fig. 5a, “illusion-corrected” in red). Importantly, γ -Net weights were optimized during training, but the orientation decoder was not. Thus, improving performance for decoding the orientation of these illusory stimuli comes at the expense of changing γ -Net weights that were responsible for its orientation-tilt illusion. As a control, another γ -Net was trained with the same routine to decode the orientation of full-image gratings, for which there is no illusion (Fig. 5a, “domain-transfer control” in blue; see Fig. S5 for training performance of both models). Both models were tested on the BSDS500 test set.

Correcting the orientation-tilt illusion of a γ -Net significantly hurts its object contour detection performance (Fig. 5a; 1-sample T -test of the per-image ODS F1 difference between models, $T(199) = 13.570, p < 0.001$). The illusion reflects γ -Net strategies for selecting object-boundaries rather than low-level contours (Fig. 5b; Fig. S6 for more examples).

6 CONCLUSION

Why do we experience visual illusions? Our experiments indicate that one representative contextual illusion, the orientation-tilt illusion, is a consequence of neural strategies for efficient scene segmentation. We directly tested the “function” of this contextual illusion with the γ -Net : a dense prediction model with recurrent modules inspired by neural circuits found in the visual cortex. The γ -Net exhibited an orientation-tilt illusion which biased it towards high-level object-boundary contours over low-level edges. On separate contour detection tasks, the γ -Net performed on par with state-of-the-art models when trained in typical “big data” regimes, but was far more efficient than these models when trained on sample-limited versions of the same datasets.

More generally, our work demonstrates novel synergy between artificial vision and vision neuroscience: we demonstrated that circuit-level insights from biology can improve the sample efficiency of deep learning models. The neural circuit that inspired the fGRU module explained biological illusions in color, motion, and depth processing (Mély et al., 2018), and we suspect that γ -Nets will have similar success in learning sample-efficient strategies – and exhibiting contextual illusions – in these domains. We will release our code upon publication to support a broader understanding of the relationship between visual illusions and neural routines for robust and efficient perception.

REFERENCES

- K. L. Briggman and D. D. Bock. Volume electron microscopy for neuronal circuit reconstruction. *Curr. Opin. Neurobiol.*, 22(1):154–161, February 2012.
- T. Coorjans, N. Ballas, C. Laurent, Ç. Gülçehre, and A. Courville. Recurrent batch normalization. In *International Conference on Learning Representations*, 2017.
- R. Deng, C. Shen, S. Liu, H. Wang, and X. Liu. Learning to predict crisp boundaries. In *Computer Vision – ECCV 2018*, pp. 570–586. Springer International Publishing, 2018.
- H. Ding, R. G. Smith, A. Poleg-Polsky, J. S. Diamond, and K. L. Briggman. Species-specific wiring for direction selectivity in the mammalian retina. *Nature*, 535(7610):105–110, July 2016.
- D. J. Felleman and D. C. Van Essen. Distributed hierarchical processing in the primate cerebral cortex. *Cereb. Cortex*, 1(1):1–47, 1991.
- A. M. Fyall, Y. El-Shamayleh, H. Choi, E. Shea-Brown, and A. Pasupathy. Dynamic representation of partially occluded objects in primate prefrontal and visual cortex. *Elife*, 6, September 2017.
- R. Geirhos, C. R. M. Temme, J. Rauber, H. H. Schütt, M. Bethge, and F. A. Wichmann. Generalisation in humans and deep neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 31*, pp. 7549–7561. Curran Associates, Inc., 2018.
- C. D. Gilbert and M. Sigman. Brain states: top-down influences in sensory processing. *Neuron*, 54(5):677–696, June 2007.
- A. Graves, S. Fernández, and J. Schmidhuber. Multi-dimensional recurrent neural networks. In *Artificial Neural Networks – ICANN 2007*, pp. 549–558. Springer Berlin Heidelberg, 2007.
- R. H. Hahnloser, R. Sarpeshkar, M. a. Mahowald, R. J. Douglas, and H. S. Seung. Digital selection and analogue amplification coexist in a cortex-inspired silicon circuit. *Nature*, 405(6789):947–951, June 2000.
- S. Hallman and C. C. Fowlkes. Oriented edge forests for boundary detection. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1732–1740, June 2015.
- J. He, S. Zhang, M. Yang, Y. Shan, and T. Huang. Bi-Directional cascade network for perceptual edge detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3828–3837, 2019.
- K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- J. C. Heck and F. M. Salem. Simplified minimal gated unit variations for recurrent neural networks. In *2017 IEEE 60th International Midwest Symposium on Circuits and Systems (MWSCAS)*, pp. 1593–1596, August 2017.
- S. Hochstein and M. Ahissar. View from the top: Hierarchies and reverse hierarchies in the visual system. *Neuron*, 36(5):791–804, 2002.
- M. Januszewski and V. Jain. Segmentation-Enhanced CycleGAN. February 2019.
- H. Kafaligonul, B. G. Breitmeyer, and H. Ögmen. Feedforward and feedback processes in vision. *Front. Psychol.*, 6:279, March 2015.
- K. Kar, J. Kubilius, K. Schmidt, E. B. Issa, and J. J. DiCarlo. Evidence that recurrent circuits are critical to the ventral stream’s execution of core object recognition behavior. *Nat. Neurosci.*, April 2019.
- N. Kasthuri, K. J. Hayworth, D. R. Berger, R. L. Schalek, J. A. Conchello, S. Knowles-Barley, D. Lee, A. Vázquez-Reina, V. Kaynig, T. R. Jones, M. Roberts, J. L. Morgan, J. C. Tapia, H. S. Seung, W. G. Roncal, J. T. Vogelstein, R. Burns, D. L. Sussman, C. E. Priebe, H. Pfister, and J. W. Lichtman. Saturated reconstruction of a volume of neocortex. *Cell*, 162(3):648–661, July 2015.

- S. W. Keemink and M. C. W. van Rossum. A unified account of tilt illusions, association fields, and contour detection based on elastica. *Vision research*, 126:164–173, September 2016.
- T. C. Kietzmann, C. J. Spoerer, L. Sorensen, and others. Recurrence required to capture the dynamic computations of the human ventral visual stream. *arXiv preprint arXiv*, 2019.
- J. K. Kim, M. Ricci, and T. Serre. Not-So-CLEVR: learning same–different relations strains feedforward neural networks. *Interface Focus theme issue on “Understanding images in biological and computer vision”*, 2018.
- J. Kim, J. K. Lee, and K. M. Lee. Deeply-Recursive convolutional network for image Super-Resolution. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1637–1645. IEEE, June 2016.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- P. C. Klink, B. Dagnino, M.-A. Gariel-Mathis, and P. R. Roelfsema. Distinct feedforward and feedback effects of microstimulation in visual cortex reveal neural mechanisms of texture segregation. *Neuron*, June 2017.
- I. Kokkinos. Pushing the boundaries of boundary detection using deep learning. *arXiv preprint arXiv:1511.07386*, 2015.
- N. Kriegeskorte. Deep neural networks: A new framework for modeling biological vision and brain information processing. *Annu Rev Vis Sci*, 1:417–446, November 2015.
- B. M. Lake, R. Salakhutdinov, and J. B. Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, December 2015.
- K. Lee, J. Zung, P. Li, V. Jain, and H. Sebastian Seung. Superhuman accuracy on the SNEMI3D connectomics challenge. In *Neural Information Processing Systems*, 2017.
- T. S. Lee and D. Mumford. Hierarchical bayesian inference in the visual cortex. *Journal of the Optical Society of America. A, Optics, image science, and vision*, 20(7):1434–1448, July 2003.
- M. Liang and X. Hu. Recurrent convolutional neural network for object recognition. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3367–3375. IEEE Computer Society, June 2015.
- D. Linsley, J. K. Kim, V. Veerabdran, C. Windolf, and T. Serre. Learning long-range spatial dependencies with horizontal gated recurrent units. In *Neural Information Processing Systems (NIPS)*, 2018a.
- D. Linsley, J. Kim, V. Veerabdran, C. Windolf, and T. Serre. Learning long-range spatial dependencies with horizontal gated recurrent units. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 31*, pp. 152–164. Curran Associates, Inc., 2018b.
- D. Linsley, D. Shiebler, S. Eberhardt, and T. Serre. Learning what and where to attend with humans in the loop. In *International Conference on Learning Representations*, 2019.
- Y. Liu, M.-M. Cheng, X. Hu, J.-W. Bian, L. Zhang, X. Bai, and J. Tang. Richer convolutional features for edge detection. *IEEE transactions on pattern analysis and machine intelligence*, 41(8): 1939–1946, August 2019.
- D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*, volume 2, pp. 416–423 vol.2, July 2001.
- D. R. Martin, C. C. Fowlkes, and J. Malik. Learning to detect natural image boundaries using local brightness, color, and texture cues. *IEEE Trans. Pattern Anal. Mach. Intell.*, 26(5):530–549, May 2004.

- D. A. Mély, D. Linsley, and T. Serre. Complementary surrounds explain diverse contextual phenomena across visual modalities. *Psychol. Rev.*, 2018.
- M. C. Mozer. Induction of multiscale temporal structure. In J. E. Moody, S. J. Hanson, and R. P. Lippmann (eds.), *Advances in Neural Information Processing Systems 4*, pp. 275–282. Morgan-Kaufmann, 1992.
- A. Nayebi, D. Bear, J. Kubilius, K. Kar, S. Ganguli, D. Sussillo, J. J. DiCarlo, and D. L. Yamins. Task-Driven convolutional recurrent models of the visual system. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 31*, pp. 5295–5306. Curran Associates, Inc., 2018.
- A. Newell, K. Yang, and J. Deng. Stacked hourglass networks for human pose estimation. In *Computer Vision – ECCV 2016*, pp. 483–499. Springer International Publishing, 2016.
- J. Nunez-Iglesias, R. Kennedy, T. Parag, J. Shi, and D. B. Chklovskii. Machine learning of hierarchical clustering to segment 2D and 3D images. *PLoS One*, 8(8):e71715, August 2013.
- R. C. O’Reilly, D. Wyatte, S. Herd, B. Mingus, and D. J. Jilk. Recurrent processing during object recognition. *Front. Psychol.*, 4(April):1–14, 2013.
- B. O’Toole and P. Wenderoth. The tilt illusion: repulsion and attraction effects in the oblique meridian. *Vision Res.*, 17(3):367–374, 1977.
- P. J. Phillips, A. N. Yates, Y. Hu, C. A. Hahn, E. Noyes, K. Jackson, J. G. Cavazos, G. Jeckeln, R. Ranjan, S. Sankaranarayanan, J.-C. Chen, C. D. Castillo, R. Chellappa, D. White, and A. J. O’Toole. Face recognition accuracy of forensic examiners, superrecognizers, and face recognition algorithms. *Proc. Natl. Acad. Sci. U. S. A.*, 115(24):6171–6176, June 2018.
- B. Recht, R. Roelofs, L. Schmidt, and V. Shankar. Do CIFAR-10 classifiers generalize to CIFAR-10? June 2018.
- P. R. Roelfsema. Cortical algorithms for perceptual grouping. *Annu. Rev. Neurosci.*, 29:203–227, January 2006.
- O. Ronneberger, P. Fischer, and T. Brox. U-Net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, pp. 234–241. Springer International Publishing, 2015.
- A. Rosenfeld, R. Zemel, and J. K. Tsotsos. The elephant in the room. August 2018.
- T. Serre. Deep learning: The good, the bad, and the ugly. *Annu Rev Vis Sci*, 5:399–426, September 2019.
- M. Siegel, T. J. Buschman, and E. K. Miller. Cortical information flow during flexible sensorimotor decisions. *Science*, 348(6241):1352–1355, June 2015.
- C. Tallec and Y. Ollivier. Can recurrent neural networks warp time? In *International Conference on Learning Representations*, 2018.
- H. Tang, M. Schrimpf, W. Lotter, C. Moerman, A. Paredes, J. Ortega Caro, W. Hardesty, D. Cox, and G. Kreiman. Recurrent computations for visual pattern completion. *Proc. Natl. Acad. Sci. U. S. A.*, 115(35):8835–8840, August 2018.
- D. Ulyanov, A. Vedaldi, and V. Lempitsky. Instance normalization: The missing ingredient for fast stylization. July 2016.
- E. Vorontsov, C. Trabelsi, S. Kadoury, and C. Pal. On orthogonality and learning recurrent networks with long term dependencies. January 2017.
- Y. Wang, X. Zhao, Y. Li, and K. Huang. Deep crisp boundaries: From boundaries to Higher-Level tasks. *IEEE transactions on image processing: a publication of the IEEE Signal Processing Society*, 28(3):1285–1298, March 2019.

- D. Wyatte, D. J. Jilk, and R. C. O'Reilly. Early recurrent feedback facilitates visual object recognition under challenging conditions. *Front. Psychol.*, 5:674, July 2014.
- S. Xie and Z. Tu. Holistically-Nested edge detection. *International journal of computer vision*, 125(1-3):3–18, December 2017.
- D. L. K. Yamins and J. J. DiCarlo. Using goal-driven deep learning models to understand sensory cortex. *Nat. Neurosci.*, 19(3):356–365, February 2016.

A γ -NET

Connectomics The current standard for computer vision applications in connectomics is to train and test on separate partitions of the same tissue volume (Januszewski & Jain, 2019). This makes it difficult to develop new model architectures without overfitting to any particular dataset. For this reason, we first tuned our connectomics γ -Net and its hyperparameters on a synthetic dataset of cell images (data not shown).

In our experiments on synthetic data, we noted monotonically improved performance with increasing timesteps, which motivated our choice of building these models with as many timesteps as could fit into GPU memory without sacrificing (we carried this concept over to the design of our γ -Nets for BSDS). Thus, we settled on 8 timesteps for the γ -Nets. We also compared our use of fGRU modules to learn recurrent connection vs. the classic LSTM and GRU recurrent modules, and found that the γ -Net was far more effective on small datasets, which we take as evidence that its recurrent application of suppression separately from facilitation is a better inductive bias for learning contour tasks (see Hahnloser et al. 2000 for a theoretical discussion on how these operations can amount to a digital selection of task-relevant features through inhibition, followed by an analog amplification of the residuals through excitation).

We found that γ -Net cell membrane detection was improved when every bottom-up unit (from a typical convolution) was given a hidden state. Like with gated recurrent architectures, these gates enable gradients to effectively skip timesteps of processing where they pathologically decay. We do this by converting every convolutional layer (except the first and last) into a “minimal gated unit” (Heck & Salem 2017). This conversion introduced two additional kernels to each convolutional layer, $U^F, W^H \in \mathbb{R}^{1 \times 1 \times K \times K}$, where the former was responsible for selecting channels from a persistent activity $\mathbf{H} \in \mathbb{R}^{X \times Y \times K}$ for processing on a given timestep and updating the persistent activity. The latter kernel transformed a modulated version of the hidden state \mathbf{H} . This transformed hidden state was combined with a vanilla convolutional feedforward encoding, $\mathbf{Z} \in \mathbb{R}^{X \times Y \times K}$ (see Eq A for the treatment). Weights in these layers were initialized with orthogonal random initializations, which help training recurrent networks (Vorontsov et al., 2017).

$$\begin{aligned} \mathbf{F} &= \sigma(\mathbf{Z} + W^F * \mathbf{H}[t - 1] + \mathbf{b}_F) \\ \mathbf{H}[t] &= \mathbf{F} \odot \mathbf{H}[t - 1] + (1 - \mathbf{F}) \odot \text{ELU}(\mathbf{Z} + W^H * (\mathbf{F} \odot \mathbf{H}[t - 1]) + \mathbf{b}_H) \end{aligned} \tag{1}$$

fGRU Here we describe additional details of the fGRU. fGRU kernels for computing suppressive and facilitative interactions have symmetric weights between channels, similar to the original circuit of Mély et al. (2018). This means that the weight $W_{x_0+\Delta x, y_0+\Delta y, k_1, k_2}$ is equal to the weight $W_{x_0+\Delta x, y_0+\Delta y, k_2, k_1}$, where x_0 and y_0 denote kernel center. This constraint means that there are nearly half as many learnable connections as a normal convolutional kernel. In our experiments, this constraint improved performance.

BSDS (20M parameters) γ -Net		
Layer	Operation	Output shape
conv-1-down	conv $3 \times 3 / 1$	$320 \times 480 \times 64$
	conv $3 \times 3 / 1$	$320 \times 480 \times 64$
	maxpool $2 \times 2 / 2$	$160 \times 240 \times 64$
conv-2-down	conv $3 \times 3 / 1$	$160 \times 240 \times 128$
	conv $3 \times 3 / 1$	$160 \times 240 \times 128$
	fGRU-horizontal $3 \times 3 / 1$	$160 \times 240 \times 128$
	maxpool $2 \times 2 / 2$	$80 \times 120 \times 128$
conv-3-down	conv $3 \times 3 / 1$	$80 \times 120 \times 256$
	conv $3 \times 3 / 1$	$80 \times 120 \times 256$
	conv $3 \times 3 / 1$	$80 \times 120 \times 256$
	fGRU-horizontal $3 \times 3 / 1$	$80 \times 120 \times 256$
	maxpool $2 \times 2 / 2$	$40 \times 60 \times 256$
conv-4-down	conv $3 \times 3 / 1$	$40 \times 60 \times 512$
	conv $3 \times 3 / 1$	$40 \times 60 \times 512$
	conv $3 \times 3 / 1$	$40 \times 60 \times 512$
	fGRU-horizontal $3 \times 3 / 1$	$40 \times 60 \times 512$
	maxpool $2 \times 2 / 2$	$20 \times 30 \times 512$
conv-5-down	conv $3 \times 3 / 1$	$20 \times 30 \times 512$
	conv $3 \times 3 / 1$	$20 \times 30 \times 512$
	conv $3 \times 3 / 1$	$20 \times 30 \times 512$
	fGRU-horizontal $3 \times 3 / 1$	$20 \times 30 \times 512$
conv-4-up	instance-norm	$20 \times 30 \times 512$
	bilinear-resize	$40 \times 60 \times 512$
	conv $1 \times 1 / 1$	$40 \times 60 \times 8$
	conv $1 \times 1 / 1$	$40 \times 60 \times 512$
	fGRU-top-down $1 \times 1 / 1$	$40 \times 60 \times 512$
conv-3-up	instance-norm	$40 \times 60 \times 512$
	bilinear-resize	$80 \times 120 \times 512$
	conv $1 \times 1 / 1$	$80 \times 120 \times 16$
	conv $1 \times 1 / 1$	$80 \times 120 \times 256$
	fGRU-top-down $1 \times 1 / 1$	$80 \times 120 \times 256$
conv-2-up	instance-norm	$80 \times 120 \times 256$
	bilinear-resize	$160 \times 240 \times 256$
	conv $1 \times 1 / 1$	$160 \times 240 \times 64$
	conv $1 \times 1 / 1$	$160 \times 240 \times 128$
	fGRU-top-down $1 \times 1 / 1$	$160 \times 240 \times 128$
Readout	instance-norm	$160 \times 240 \times 128$
	bilinear-resize	$320 \times 480 \times 128$
	conv $1 \times 1 / 1$	$320 \times 480 \times 1$

Table S1: γ -Net architecture for contour detection in BSDS natural images. Down refers to down-sampling layers; up refers to up-sampling layers, and readout maps model activities into per-pixel decisions. Kernels are described as kernel-height \times kernel-width / stride size. All convolutional layers except for the Readout use non-linearities. All non-linearities in this network are linear rectifications. Model predictions come from the fGRU hidden state for conv-2-down, which are resized to match the input image resolution and passed to the linear per-pixel readout.

Connectomics γ -Net (450K parameters)		
Layer	Operation	Output shape
conv-1-down	conv $3 \times 3 / 1$	$384 \times 384 \times 24$
	conv $3 \times 3 / 1$	$384 \times 384 \times 24$
	fGRU-horizontal $9 \times 9 / 1$	$384 \times 384 \times 24$
	maxpool $2 \times 2 / 2$	$192 \times 192 \times 24$
conv-2-down	conv $3 \times 3 / 1$	$192 \times 192 \times 28$
	fGRU-horizontal $7 \times 7 / 1$	$192 \times 192 \times 28$
	maxpool $2 \times 2 / 2$	$96 \times 96 \times 28$
conv-3-down	conv $3 \times 3 / 1$	$96 \times 96 \times 36$
	fGRU-horizontal $5 \times 5 / 1$	$96 \times 96 \times 36$
	maxpool $2 \times 2 / 2$	$48 \times 48 \times 36$
conv-4-down	conv $3 \times 3 / 1$	$48 \times 48 \times 48$
	fGRU-horizontal $3 \times 3 / 1$	$48 \times 48 \times 48$
	maxpool $2 \times 2 / 2$	$24 \times 24 \times 48$
conv-5-down	conv $3 \times 3 / 1$	$24 \times 24 \times 64$
	fGRU-horizontal $1 \times 1 / 1$	$24 \times 24 \times 64$
conv-4-up	transpose-conv $4 \times 4 / 2$	$48 \times 48 \times 48$
	conv $3 \times 3 / 1$	$48 \times 48 \times 48$
	instance-norm	$48 \times 48 \times 48$
	fGRU-top-down $1 \times 1 / 1$	$48 \times 48 \times 48$
conv-3-up	transpose-conv $4 \times 4 / 2$	$96 \times 96 \times 36$
	conv $3 \times 3 / 1$	$96 \times 96 \times 36$
	instance-norm	$96 \times 96 \times 36$
	fGRU-top-down $1 \times 1 / 1$	$96 \times 96 \times 36$
conv-2-up	transpose-conv $4 \times 4 / 2$	$192 \times 192 \times 28$
	conv $3 \times 3 / 1$	$192 \times 192 \times 28$
	instance-norm	$192 \times 192 \times 28$
	fGRU-top-down $1 \times 1 / 1$	$192 \times 192 \times 28$
conv-1-up	transpose-conv $4 \times 4 / 2$	$384 \times 384 \times 24$
	conv $3 \times 3 / 1$	$384 \times 384 \times 24$
	instance-norm	$384 \times 384 \times 24$
	fGRU-top-down $1 \times 1 / 1$	$384 \times 384 \times 24$
Readout	instance-norm	$384 \times 384 \times 24$
	conv $5 \times 5 / 1$	$384 \times 384 \times 24$

Table S2: γ -Net architecture for cell membrane detection in SEM images. Down refers to down-sampling layers; up refers to up-sampling layers, and readout maps model activities into per-pixel decisions. Kernels are described as kernel-height \times kernel-width / stride size. All fGRU non-linearities are linear rectifications, and all convolutional non-linearities are exponential linear units (ELU), as in (Lee et al., 2017). All convolutional layers except for the Readout use non-linearities. Model predictions come from the fGRU hidden state for conv-1-down, which are passed to the linear readout.

While optimizing γ -Nets on synthetic cell image datasets, we found that a small modification of the fGRU input gate offered a modest improvement in performance. We realized that the input gate in the fGRU is conceptually similar to recently developed models for feedforward self-attention in deep neural networks. Specifically, the global-and-local attention modules of (Linsley et al., 2019), in which a non-linear transformation of a layer’s activity is used to modulate the original activity. Here, we took inspiration from global-and-local attention, and introduced an additional gate into the fGRU, resulting in the following modification of the main equations.

Stage 1:

$$\begin{aligned}
 \mathbf{A}^S &= U^A * \mathbf{H}[t-1] && \# \text{ Compute channel-wise selection} \\
 \mathbf{M}^S &= U^M * \mathbf{H}[t-1] && \# \text{ Compute spatial selection} \\
 \mathbf{G}^S &= \text{sigmoid}(IN(\mathbf{A}^S \odot \mathbf{M}^{S*})) && \# \text{ Compute suppression gate} \\
 \mathbf{C}^S &= IN(W^S * (\mathbf{H}[t-1] \odot \mathbf{G}^S)) && \# \text{ Compute suppression interactions} \\
 \mathbf{S} &= \left[\mathbf{Z} - \left[(\alpha \mathbf{H}[t-1] + \mu) \mathbf{C}^S \right]_+ \right]_+, && \# \text{ Additive and multiplicative suppression of } \mathbf{Z}
 \end{aligned}$$

Stage 2:

$$\begin{aligned}
 \mathbf{G}^F &= \text{sigmoid}(IN(U^F * \mathbf{S})) && \# \text{ Compute channel-wise recurrent updates} \\
 \mathbf{C}^F &= IN(W^F * \mathbf{S}) && \# \text{ Compute facilitation interactions} \\
 \tilde{\mathbf{H}} &= \left[\nu(\mathbf{C}^F + \mathbf{S}) + \omega(\mathbf{C}^F * \mathbf{S}) \right]_+ && \# \text{ Additive and multiplicative facilitation of } \mathbf{S} \\
 \mathbf{H}[t] &= (1 - \mathbf{G}^F) \odot \mathbf{H}[t-1] + \mathbf{G}^F \odot \tilde{\mathbf{H}} && \# \text{ Update recurrent state}
 \end{aligned}$$

where $IN(\mathbf{r}; \delta, \nu) = \nu + \delta \odot \frac{\mathbf{r} - \hat{\mathbb{E}}[\mathbf{r}]}{\sqrt{\widehat{\text{Var}}[\mathbf{r}] + \eta}}$.

This yields the global input gate activity $\mathbf{A}^S \in \mathbb{R}^{X \times Y \times K}$ and the local input gate activity $\mathbf{M}^{S*} \in \mathbb{R}^{X \times Y \times I}$, which are computed as filter responses between the previous hidden state $\mathbf{H}[t-1]$ and the global gate kernel $U^A \in \mathbb{R}^{1 \times 1 \times K \times K}$ and the local gate kernel $U^M \in \mathbb{R}^{3 \times 3 \times K \times I}$. Note that the latter filter is learning a mapping into 1 dimension and is therefore first tiled into K dimensions, yielding \mathbf{M}^{S*} , before elementwise multiplication with \mathbf{A}^S . All results in the main text use this implementation.

Following the lead of (Linsley et al., 2018a), we incorporated normalizations into the fGRU. Let $\mathbf{r} \in \mathbb{R}^d$ denote the vector of layer activations that will be normalized. We chose instance normalization (Ulyanov et al., 2016) since it is independent of batch size, which was 1 for γ -Nets in our experiments. Instance normalization introduces two k -dimensional learned parameters, $\delta, \nu \in \mathbb{R}^d$, which control the scale and bias of normalized activities, and are shared across timesteps of processing. In contrast, means and variances are computed on every timestep, since fGRU activities are not i.i.d. across timesteps. Elementwise multiplication is denoted by \odot and η is a regularization hyperparameter.

Learnable gates, such as those in the fGRU, are helpful for training RNNs. But there are other heuristics that are also important for optimizing performance. We use several of these with γ -Nets, such as Chronos initialization of fGRU gate biases (Tallec & Ollivier, 2018) and random orthogonal initialization of kernels (Vorontsov et al., 2017). We initialized the learnable scale parameter δ of fGRU normalizations to 0.1, since values near 0 optimize the dynamic range of gradients passing through its sigmoidal gates (Cooijmans et al., 2017). Similarly, fGRU parameters for learning additive suppression/facilitation (μ, ν) were initialized to 0, and parameters for learning multiplicative inhibition/excitation (α, ω) were initialized to 0.1. Finally, when implementing top-down connections, we incorporated an extra skip connection. The activity of layer ℓ was added to the fGRU-computed top-down interactions between layer ℓ and layer $\ell + 1$. This additional skip connection improved the stability of training.

Name	Tissue	Imaging	Resolution	Voxels (X/Y/Z/Volumes)
SNEMI3D	Mouse cortex	mbSEM	$6 \times 6 \times 29\text{nm}$	$1024 \times 1024 \times 100 \times 1$
Ding	Mouse retina	SBEM	$13.2 \times 13.2 \times 26\text{nm}$	$384 \times 384 \times 384 \times 1$

Table S3: SEM image volumes used in membrane prediction. SNEMI3D images and annotations are publicly available (Kashuri et al., 2015), whereas the Ding dataset is a volume from (Ding et al., 2016) that we annotated.

B MEMBRANE PREDICTION MODELS

Our reference model for membrane prediction is the 3D U-Net of (Lee et al., 2017). We followed their published routine for validating our implementation of their U-Net model.

Key to the approach of Lee et al. (2017) is their use of a large set of random data augmentations applied to SEM image volumes, which simulate common noise and errors in SEM imaging. These are (i) misalignment between consecutive z -locations in each input image volume. (ii) Partial- or fully-missing sections of the input image volumes. (iii) Blurring of portions of the image volume. Augmentations that simulated these types of noise, as well as random flips over the xyz -plane, rotations by 90° , brightness and contrast perturbations, were applied to volumes following the settings of Lee et al. (2017). The model was trained using Adam (Kingma & Ba, 2014) and the learning rate schedule of Lee et al. (2017), in which the optimizer step-size was halved when validation loss stopped decreasing (up to four times). Training involved single-SEM volume batches of $160 \times 160 \times 18$ (X/Y/Z), normalized to $[0, 1]$. As in Lee et al. (2017), models were trained to predict nearest-neighbor voxel affinities, as well as 3 other mid- to long-range voxel distances. Only nearest neighbor affinities were used at test time.

C ORIENTATION-TILT ILLUSION IMAGE DATASET

Models were tested for a tilt illusion by first training on grating images of a single orientation, then testing on images in which a center grating had the same/different orientation as a surround grating. Each image in the training dataset consisted of a circular patch of oriented grating on a gray canvas of size 500×500 pixels. To ensure that the decoder successfully decoded orientation information from model activities, the training dataset incorporated a wide variety of grating stimuli with 4 randomly sampled image parameters: r , λ , θ , and ϕ . r denotes the radius of the circle in which oriented grating has been rendered, and was sampled from a uniform distribution with interval between 80 and 240 pixels; λ specifies the wavelength of the grating pattern and was sampled from a uniform distribution with interval between 30 and 90 pixels; θ specifies the orientation of the gratings and is uniformly sampled from all possible orientations; ϕ denotes the phase offset of the oriented gratings and is also uniformly sampled from all possible values. The models' BSDS-trained weights were fixed and readout layers were trained to decode orientation at the center of each image (procedure described in the main text).

This setup allowed us to tease apart the effects of the surround on the representation of orientation in the center by introducing separate surround regions in each test image filled with gratings with same/different orientations as the center (Fig.S2b). Each test image was generated with one additional parameter, $\Delta\theta$ which specified orientation difference of the surround gratings with respect to the center orientation, θ , and was sampled from a uniform distribution with interval between -90 and $+90$ degrees. The radius of the surround grating is denoted by r and was sampled from the same uniform distribution we used in training dataset. Center gratings are then rendered in a circle of radius that is one half of the surround gratings.

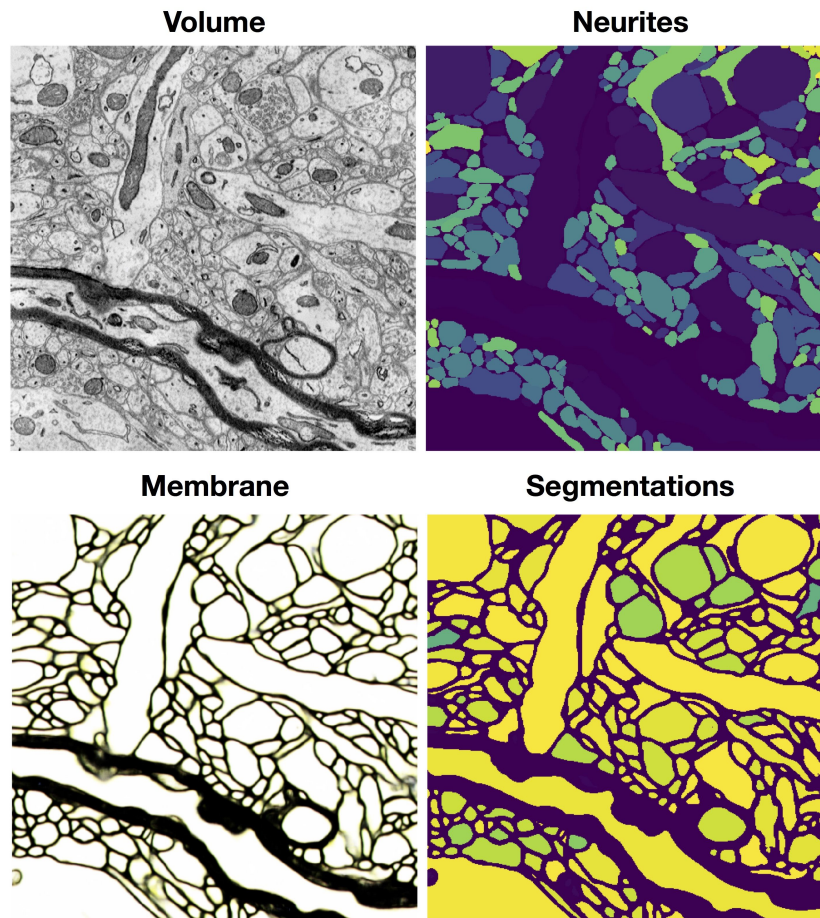


Figure S1: We trained the reference 3D U-Net from (Lee et al., 2017) on the SNEMI3D dataset to validate the implementation. Segmentations here are derived by watershedding and agglomeration with GALA (Nunez-Iglesias et al., 2013), resulting in “superhuman” ARAND (evaluated according to the SNEMI3D standard; lower is better) of 0.04, which is below the reported human-performance threshold of 0.06 and on par with the published result (see Table 1 in Lee et al. 2017, mean affinity agglomeration).

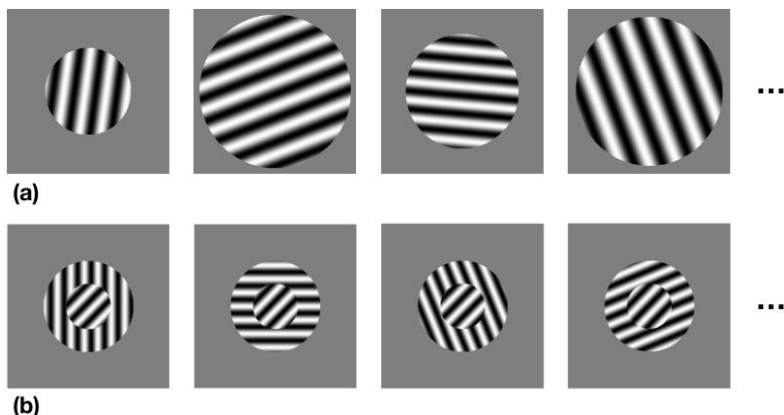


Figure S2: Examples of tilt-illusion stimuli. **(a)** For training images, we sample over a range of size and wavelength to generate single oriented grating patches. **(b)** Test images are obtained by sampling a full range of surround orientation, while fixing all other parameters such as size and frequency of gratings as well as the orientation of the center gratings (at 45 degrees).

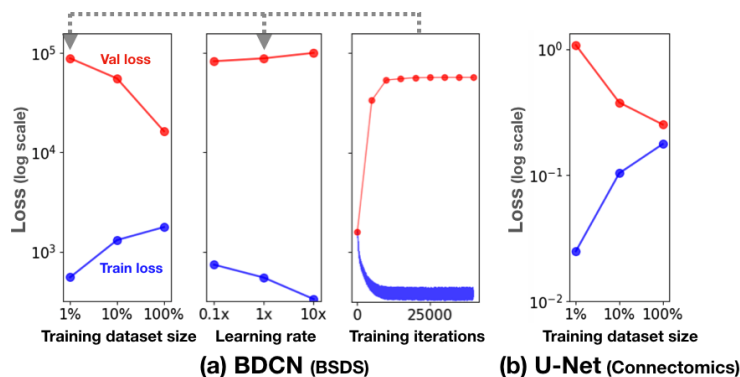


Figure S3: Searching over learning rates did not rescue BSDS performance on small BSDS datasets had no affect on performance. over learning rates did not rescue from overfitting. The left panel depicts training and validation losses for BSDS on different sized subsets of BSDS500. **(b)** Performance after training with three different learning rates on the 5% split. There is little difference in best validation performance between the three learning rates. **(c)** The full training and validation loss curves for the BDCN trained on 5% of BSDS. The model overfits immediately. The model also overfit on the other dataset sizes, but because there was more data, this happened later in training.

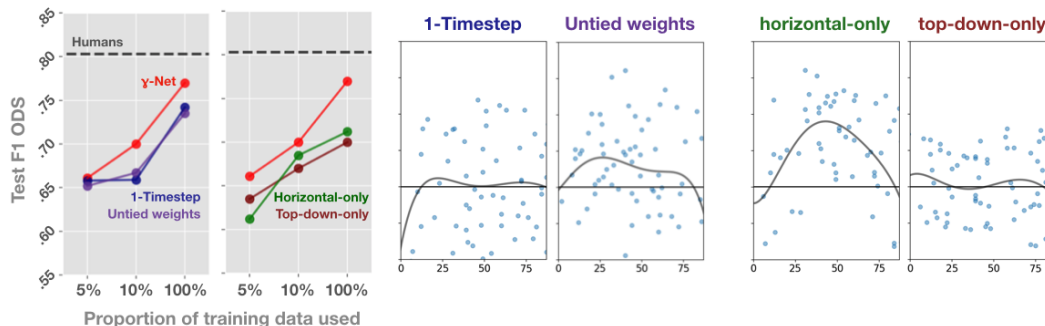


Figure S4: Performance of lesioned γ -Nets on BSDS and tests for orientation-tilt illusions. The first BSDS performance panel depicts performance of “feedforward” lesioned versions of γ -Net. These are a γ -Net run for 1 timestep, and a γ -Net which had unique/“untied” weights on each of its 8-timesteps of processing (similar to the stacked hourglass network of Newell et al. 2016). The second BSDS performance panel depicts “recurrent” lesioned versions of γ -Net. These are a version that can only learn horizontal connections (no top-down fGRUs), and a version that can only learn top-down connections (all fGRUs used 1×1 spatial kernels). Also shown are tests for the orientation-tilt illusion in each of these models. The horizontal-only model shows a relatively apparent repulsion regime, but does not also exhibit the attraction regime that defines the illusion.

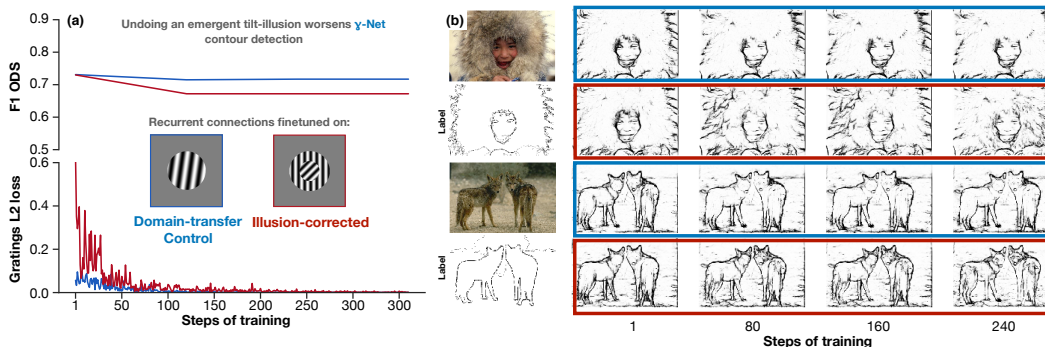


Figure S5: Performance of γ -Nets during experiments to correct an orientation-tilt illusion. The illusion-corrected model was trained to have veridical representations of the central grating in tilt-illusion stimuli. To control for potential detrimental effects of the training procedure per se, a control model (“domain-transfer control”) was trained to decode orientations of single-grating stimuli. (a) Training causes contour-detection performance of both models to drop. However, the illusion-corrected model performance drops significantly more than the biased model (see main text for hypothesis testing). The losses for both models converge towards 0 across training, indicating that both learned to decode central-orientations of their stimuli. (b) Contour detection examples for biased and bias-corrected models across steps of this training procedure.

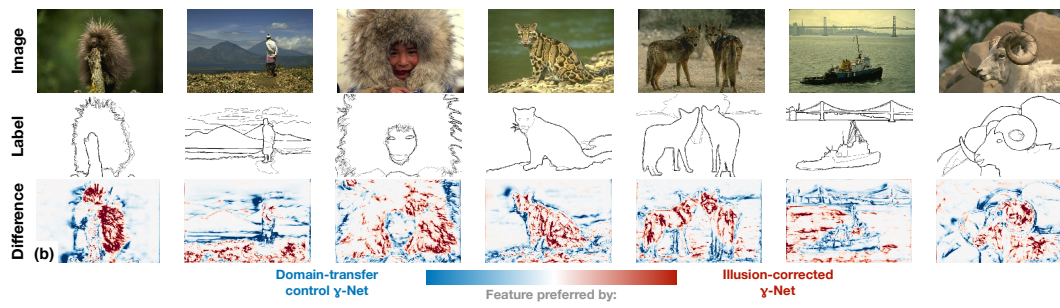


Figure S6: Differences in contour predictions for the illusion-corrected and domain-transfer control γ -Nets on BSDS500.