

Capsule Networks but not Classic CNNs Explain Global Visual Processing

Adrien Doerig^{a,†,*}, Lynn Schmittwilken^{a,†}, Bilge Sayim^{b,c}, Mauro Manassi^d & Michael H. Herzog^a

^a Laboratory of Psychophysics, Brain Mind Institute, EPFL, Lausanne, 1015, Switzerland

^b Institute of Psychology, University of Bern, 3012 Bern, Switzerland

^c Univ. Lille, CNRS, UMR 9193- SCALab- Sciences Cognitives et Sciences Affectives, F-59000 Lille, France

^d Department of Psychology, University of Aberdeen, Aberdeen, Scotland, UK

[†] Equal contributions

* Corresponding author: adrien.doerig@gmail.com

Keywords: Vision, Neural Networks, Capsule Networks, Crowding, Global Shape Processing, Recurrent Processing

Abstract

Classically, visual processing is described as a cascade of local feedforward computations. Feedforward Convolutional Neural Networks (ffCNNs) have shown how powerful such models can be and revolutionized computer vision. However, ffCNNs only roughly mimic human vision. They lack recurrent connections and rely mainly on local features, contrary to humans who use global shape computations. Previously, using visual crowding as a well-controlled challenge, we showed that no classic model of vision, including ffCNNs, can explain human global shape processing (1). Here, we show that Capsule Neural Networks (CapsNets; 2), combining ffCNNs with a grouping and segmentation mechanism, solve this challenge in a natural manner. We hypothesize that one computational function of recurrence is to efficiently implement grouping and segmentation. We provide psychophysical evidence that, indeed, time-consuming recurrent processes implement complex grouping and segmentation in humans. CapsNets reproduce these results in a natural manner. Together, we provide mutually reinforcing psychophysical and computational evidence that a recurrent grouping and segmentation process is essential to understand the visual system and create better models that harness global shape computations.

28 Introduction

29 The visual system is often seen as a hierarchy of local feedforward computations (3), going back to
30 the seminal work of Hubel and Wiesel (4). Low-level neurons detect basic features, such as edges.
31 Higher-level neurons pool the outputs from the lower-level neurons to detect higher-level features
32 such as corners, shapes, and ultimately objects. Feedforward Convolutional Neural Networks (ffCNNs)
33 embody this classic framework of vision and excel at object detection (5). However, despite their
34 amazing success, ffCNNs only roughly mimic human vision. For example, they lack the abundant
35 recurrent processing of humans (6, 7), perform differently than humans in crucial psychophysical tasks
36 (1, 8), and can be easily misled (9–11). Importantly, ffCNNs focus mainly on local, texture-like features,
37 while humans harness global shape level computations (1, 11–15).

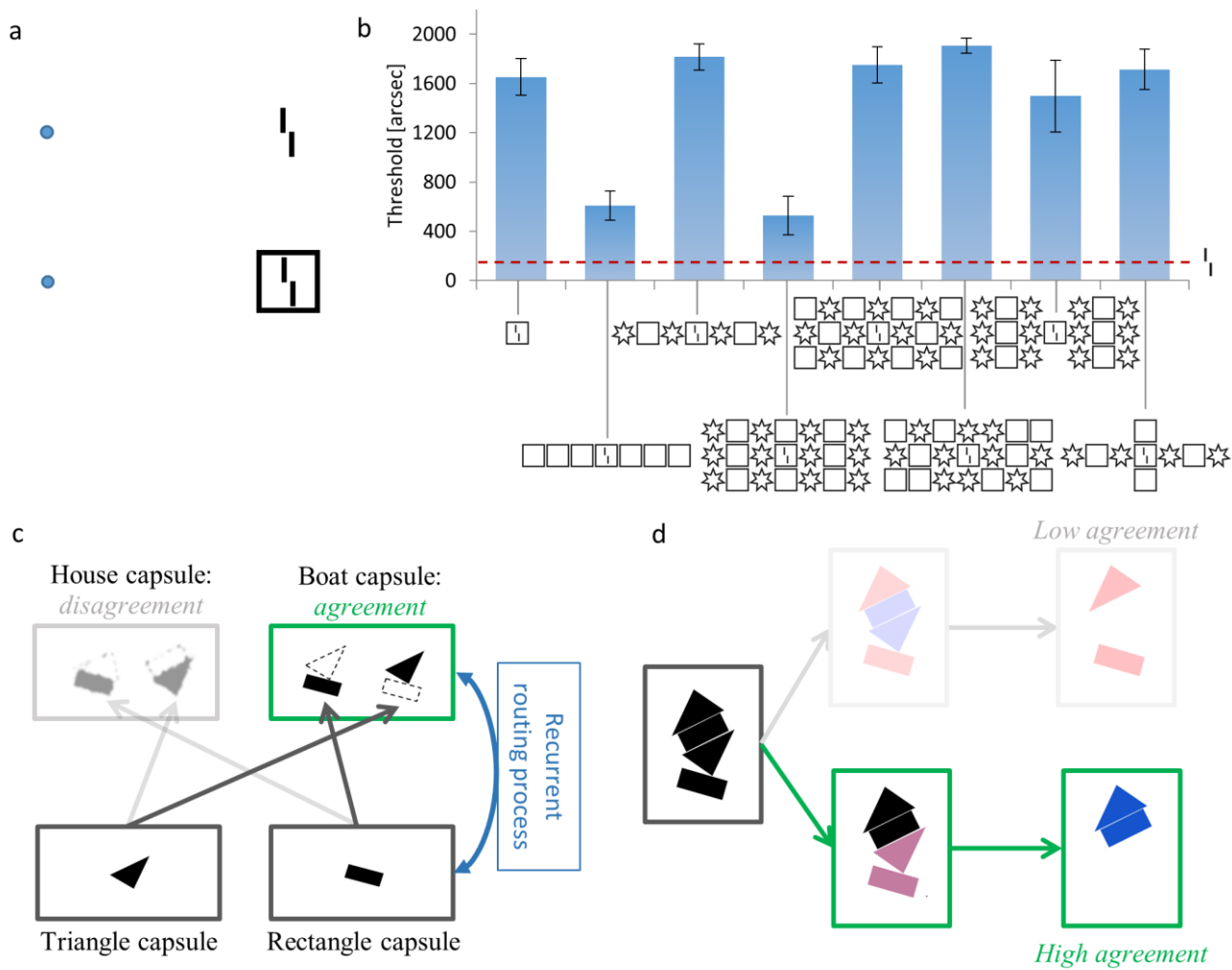
38 One difficulty in addressing these topics is that there are no widely accepted diagnostic tools to
39 specifically characterize global shape level computations in neural networks. Models are usually
40 compared either on computer vision benchmarks, such as ImageNet (16), or with neural responses in
41 the visual system (17, 18). One drawback with these approaches is that the datasets are hard to
42 control. Psychophysical results can be used to fill this gap and create well-controlled challenges for
43 visual models, tailored to target specific aspects of vision (19). Here, we use visual crowding to target
44 global shape computations in humans and machines.

45 In crowding, objects that are easy to identify in isolation seem jumbled and indistinct when clutter is
46 added (1, 20–25). For example, a vernier target is presented, i.e., two vertical lines separated by a
47 horizontal offset (Figure 1a). When the vernier is presented alone, observers easily discriminate the
48 offset direction. When a flanking square surrounds the target, performance drops, i.e., there is strong
49 crowding (26, 27). Surprisingly, *adding* more flanking squares *reduces* crowding strongly, depending
50 on the configuration (Figure 1b; 25). This global, configurational *uncrowding* effect occurs for a wide
51 range of stimuli in vision, including foveal and peripheral vision, audition, and haptics (28–34). The
52 ubiquity of (un)crowding in perception is not surprising since elements are rarely seen in isolation.
53 Hence, any perceptual system needs to cope with crowding, i.e., isolating important information from
54 clutter.

55 We have shown previously that these global effects of crowding *cannot* be explained by models based
56 on the classic framework of vision, including ffCNNs (1, 15, 35). Here, we propose a new framework
57 to understand these global computations. We show that Capsule Neural Networks (CapsNets; 2),

58 augmenting fFCNNs with a recurrent grouping and segmentation process, can explain these complex
 59 global (un)crowding results in a natural manner. Two processing regimes can occur in CapsNets: a fast
 60 feedforward pass able to quickly process information, and a time-consuming recurrent regime to
 61 perform more in depth global grouping and segmentation computations. We will show that the
 62 human visual system indeed harnesses recurrent processing for efficient grouping and segmentation,
 63 and that CapsNets naturally explain this result. Together, our results suggest that a time-consuming
 64 recurrent grouping and segmentation process is crucial for shape-level computations in both humans
 65 and artificial neural networks.

66



67

68 **Figure 1: a. Crowding:** Perception of visual elements deteriorates in clutter, an effect called crowding. In this example, a
 69 vernier (two vertical bars with a horizontal offset) becomes harder to perceive when a square flanker is added (fixate on
 70 the blue dots). **b. Uncrowding:** A vernier is presented in the visual periphery. The offset direction is easily reported (dashed
 71 red line; the y-axis shows the threshold, i.e., the minimal offset size at which observers can report the offset direction with
 72 75% accuracy). When a square flanker surrounds the vernier, performance deteriorates- a classic crowding effect. When
 73 more squares are added, performance recovers (uncrowding). Critically, the uncrowding effect depends on the global

stimulus configuration. For example, if some squares are replaced by stars, performance deteriorates again (3rd bar; 25).

c. Routing by agreement in CapsNets: Information propagates between layers of capsules through a recurrent routing process aiming to maximize agreement between capsules. Each capsule is a group of neurons whose activity vector represents the pose (such as position, orientation, etc.) of the feature it detects. In this toy example, lower-level capsules detect simple shapes such as triangles and rectangles. In the next layer, capsules have learnt combinations of these shapes. Here, the triangle capsule detects a tilted triangle and the rectangle capsule detects a tilted rectangle. Each of these capsules predicts what is represented at the next layer. For example, the triangle capsule predicts an upside-down house or a tilted boat, while the rectangle capsule predicts a tilted house or a tilted boat. The recurrent routing by agreement process routes information between the layers so that agreement is maximized. In this case, capsules agree about the tilted boat, but disagree about the house orientation. Hence, the routing by agreement suppresses activity in the house capsule and boosts activity in the boat capsule.

d. Grouping and segmentation in CapsNets: This recurrent routing by agreement process endows CapsNets with natural grouping and segmentation capabilities. Here, an ambiguous stimulus, which can be seen either as an upside-down house (top) or a house on a boat (bottom), is presented. The upside-down house interpretation leaves parts of the image unexplained and this causes disagreement. Hence, the routing by agreement will select the latter interpretations because it is the best explanation of the input and therefore maximizes agreement. Thereby, the house and boat are each grouped as an object and segmented into the corresponding higher-level capsules.

Results

Experiment 1: Crowding and Uncrowding Naturally Occur in CapsNets

In CapsNets, early convolutional layers extract basic visual features. Recurrent processing combines these features into groups and segments objects by a process called *routing by agreement*¹. The entire network is trained end-to-end through backpropagation. *Capsules* are groups of neurons representing visual features and are crucial for the routing by agreement process. Low-level capsules iteratively predict the activity of high-level capsules in a recurrent loop. If the predictions agree, the corresponding high-level capsule is activated. For example, if a capsule responds to a triangle above a rectangle detected by another capsule, they agree that the higher-level object should be a house and, therefore, the corresponding high-level capsule is activated (Figure 1c). This process allows CapsNets to group and segment objects (Figure 1d).

We trained CapsNets with two convolutional layers followed by two capsule layers to recognize greyscale images of vernier targets and groups of identical shapes (see Methods). During training, either a vernier or a group of identical shapes was presented. The network had to simultaneously

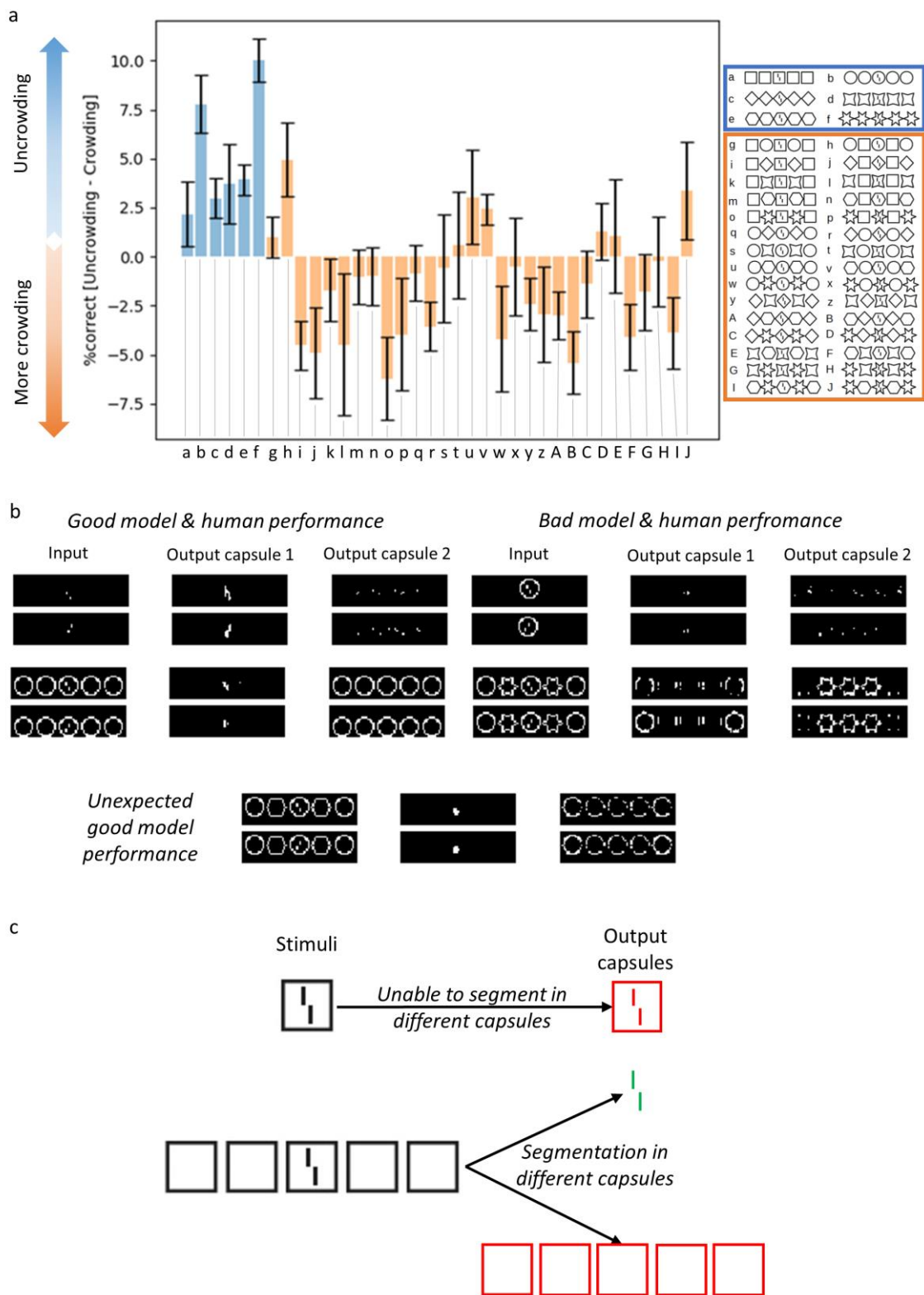
¹ In most implementations of CapsNets, including ours and (2), the iterative routing by agreement process is not explicitly implemented as a “standard” recurrent neural network processing sequences of inputs online. Instead, there is an iterative algorithmic loop (see (2) for the algorithm), which is equivalent to recurrent processing.

106 classify the shape type, the number of shapes in the group, and the vernier offset direction.
107 Importantly, verniers and shapes were never presented together during training, i.e., there were no
108 (un)crowding stimuli during training.

109 When combining verniers and shapes after training, both crowding and uncrowding occurred (Figure
110 2a): presenting the vernier target within a single flanker deteriorated vernier offset discrimination
111 (crowding), and adding more identical flankers recovered performance (uncrowding). Adding config-
112 urations of alternating different flankers did not recover the network’s performance, similarly to hu-
113 man vision. Small changes in the network hyperparameters or stimulus characteristics do not affect
114 these results (supplementary material). As a control condition, we checked that when the vernier
115 target is presented outside the flanker configuration, rather than inside, there was no performance
116 drop (supplementary material). Hence, the performance drop in crowded conditions was due to
117 crowding and not merely to the simultaneous presence of the target and flanking shape in the stim-
118 ulus.

119 Reconstructing the input image based on the network’s output (see Methods) shows that (un)crowd-
120 ing occurs through grouping and segmentation (figure 2b). Crowding occurs when the target and
121 flankers cannot be segmented and are therefore routed to the same capsule. In this case, they inter-
122 fere because a single capsule cannot represent well two objects simultaneously due to limited neural
123 resources. This mechanism is similar to pooling: information about the target is pooled with infor-
124 mation about the flankers, leading to poorer representations. However, if the flankers are segmented
125 away and represented in a different capsule, the target is released from the flankers’ deleterious ef-
126 fects and *uncrowding* occurs (Figure 2c). This segmentation can only happen if the network has learnt
127 to group the flankers into a single higher-level object represented in a different capsule than the ver-
128 nier target. Segmentation is facilitated when more flankers are added because more low-level cap-
129 sules agree about the presence of the flanker group.

130 Alternating configurations of different flankers, as in the third configuration of Figure 1b, usually do
131 not lead to uncrowding (25). In some rare cases, the network produced uncrowding with such config-
132 urations (stimuli h, u, v & J; Figure 2). Reconstructions show that in these cases the network simply
133 could not differentiate between different shapes of the flankers (e.g. between circles and hexagons),
134 and the flankers were segmented away from the target (Figure 2b). This further reinforces the notion
135 that grouping and segmentation differentiate crowding from uncrowding: whenever the network
136 reaches the conclusion that flankers form a group, segmentation is facilitated. When this happens,
137 the vernier and flankers are represented in different capsules, leading to good performance.



140 **Figure 2: a. CapsNets explain both crowding and uncrowding:** The x-axis shows the various stimuli. Performance is shown
141 on the y-axis as the % correct for each stimulus *minus* the % correct with only the central single flanker. For example, in
142 column *a*, vernier offset direction is easier to read out with 5 square flankers than with 1 square flanker, as expected. Error
143 bars are the standard error over 10 network trainings (we used 10 networks to match the typical number of observers in
144 human experiments; 25, 36). The blue bars represent configurations for which *uncrowding* is expected (blue bars larger

than 0.0 are in accordance with the human data) and orange bars represent configurations for which crowding is expected (orange bars smaller than or around 0.0 are in accordance with the human data). **b. Reconstructions:** We reconstructed the input image based on the output capsules' activities (see Methods). The reconstructions based on the two most activated capsules are shown. When the vernier is presented alone (top left), the reconstructions are good. When a single flanker is added (top right), the vernier reconstruction deteriorates (crowding) because the vernier is not well segmented from the flanker. When identical flankers are added (bottom left), the vernier reconstruction recovers, i.e., it is well segmented from the flankers (uncrowding). With different flankers (bottom right), the vernier is not represented at all in the two winning capsules (crowding). Interestingly, when the network produces "unexpected" uncrowding (i.e., the network shows uncrowding contrary to humans; bottom left), the reconstructions strongly resemble the case of "normal" uncrowding (compare middle and bottom left panels). In this case, the network was unable to notice the difference between circles and hexagons, and treated both stimuli in the same way. **c. Segmentation and (un)crowding in CapsNets:** If CapsNets can segment the vernier target away from the flankers during the recurrent routing by agreement process, uncrowding occurs. Segmentation is difficult when a single flanker surrounds the target because capsules disagree about what is shown at this location. In the case of configurations that the network has learned to group, many primary capsules agree about the presence of a group of shapes, which can therefore easily be segmented away from the vernier target.

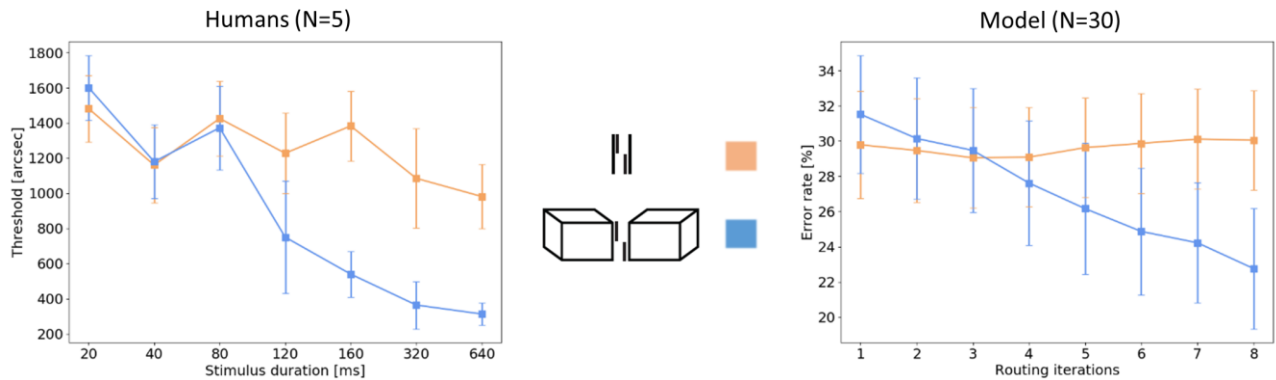
Experiment 2: The role of recurrent processing

As mentioned, processing in CapsNets starts with a feedforward sweep followed by recurrent routing by agreement to refine grouping and segmentation. We hypothesize that humans may use recurrent processing to efficiently implement grouping and segmentation. To test this hypothesis, we psychophysically investigated the temporal dynamics of (un)crowding. We show that uncrowding is mediated by a time-consuming *recurrent* process in humans. When the target groups with the flankers, crowding occurs immediately. In contrast, when the target and flankers form separate groups, time-consuming recurrent computations are required to segment the flanker from the target. We successfully model these results with CapsNets.

First, we performed a psychophysical crowding experiment with a vernier target flanked by either two lines or two cuboids (see Methods; Figure 3). The stimuli were displayed for varying durations from 20 to 640ms and five observers reported the vernier offset direction. For short stimulus durations, crowding occurred for both flanker types, i.e., thresholds increased for both the lines and cuboids conditions compared to the vernier alone condition (lines: $p = 0.0017$, cuboids: $p = 0.0013$, 2-tailed one-sample t-tests).

We quantified how performance changed with increasing stimulus duration by fitting a line $y = ax + b$ to the data for each subject, and comparing the mean slope a across subjects with 0 in one-sample 2-tailed t-tests. The performance on the lines condition did not significantly change with increasing

179 stimulus duration ($p = 0.057$). These results are in accordance with previous results which show that
 180 crowding varies very little with stimulus duration (37; but see 38, 39). With the flanking cuboids we
 181 found a different pattern of results: performance dramatically improves with stimulus duration ($p =$
 182 0.0007). This improvement cannot be explained by local mechanisms, such as lateral inhibition (26,
 183 40) or pooling (41–43) since the inner flanking vertical lines are the same in the lines and cuboids.
 184 Hence, according to a local approach we should expect no difference in thresholds between the two
 185 flanking conditions.



186
 187 **Figure 3: Temporal dynamics of uncrowding:** *Left: Human data.* For cuboid flankers, strong crowding occurs up to 100ms
 188 of stimulus presentation, and then uncrowding gradually occurs for longer durations (i.e., performance improves; blue).
 189 The x-axis shows different stimulus durations and the y-axis shows the corresponding thresholds (i.e., lower values indicate
 190 better performance). Error bars indicate standard error. Uncrowding does not occur with single line flankers, even for long
 191 stimulus durations (orange). We hypothesize that the cuboids are segmented from the vernier target through time-con-
 192 suming recurrent processing (the line flankers are grouped with the target and cannot be segmented at all). *Right: Model*
 193 *data.* CapsNets can explain these results by varying the number of recurrent routing by agreement iterations. The x-axis
 194 shows different numbers of routing iterations during testing and the y-axis shows the corresponding error rates (i.e., lower
 195 values indicate better performance). Error bars indicate standard deviation across 30 trained networks (see Methods).
 196 Similarly to humans, both lines and cuboids lead to crowding with few routing by agreement iterations. Performance
 197 increases with routing iterations only for the cuboids. This suggests that recurrent processing helps to compute and seg-
 198 ment the complex cuboids, but the lines are immediately strongly grouped with the vernier and can never be segmented.
 199 Hence, they do not benefit from the recurrent segmentation process.

200
 201 Crucially, uncrowding occurred for the cuboid flankers only when stimulus durations were sufficiently
 202 long (Figure 3). In contrast, the effect of the line flankers does not change over time. We propose that
 203 these results reflect the time-consuming recurrent computations needed to segment the cuboid
 204 flankers away from the target. Performance does not improve with the line flankers, because they are
 205 too strongly grouped with the vernier target, so recurrent processing cannot segment them away.

206 We trained CapsNets with the same architecture as in experiment 1 to discriminate vernier offsets,
207 and to recognize lines, cuboids and scrambled cuboids (see Methods; the scrambled cuboids were
208 included only to prevent the network from classifying lines vs. cuboids simply based on the number
209 of pixels in the image). As in experiment 1, during training, each training sample contained one of the
210 shape types, and the network had to classify which shape type was present and to discriminate the
211 vernier offset direction. We used 8 routing by agreement iterations during training. As in experiment
212 1, verniers and flankers were never presented together during training (i.e., there were no
213 (un)crowding stimuli).

214 After training, we tested the networks on (un)crowding stimuli, changing the number recurrent rout-
215 ing by agreement iterations from one (leading to a purely feedforward regime) to 8 iterations (a highly
216 recurrent regime; Figure 3). We found that CapsNets naturally explain the human results. Using the
217 same statistical analysis as for humans, we found that with more iterations, the cuboids are better
218 segmented from the target, and performance improves ($p = 0.003$). On the other hand, the effect of
219 the line flankers does not change over time ($p = 0.64$). These results were not affected by small
220 changes in network hyperparameters (supplementary material).

221 These findings are explained by the recurrent routing by agreement process. With cuboids, capsules
222 across an extended spatial region need to agree about the presence of a cuboid, which is then seg-
223 mented into its own capsule. This complex process requires several recurrent iterations of the routing
224 by agreement process. On the other hand, the lines are immediately strongly grouped with the vernier,
225 so further iterations of routing by agreement do not achieve successful segmentation and, hence,
226 cannot improve performance.

227

228 Discussion

229 Our results provide strong evidence that time-consuming recurrent grouping and segmentation is
230 crucial for shape-level computations in both humans and artificial neural networks. We used
231 (un)crowding as a psychophysical probe to investigate how the brain flexibly forms object
232 representations. These results specifically target global, shape-level and time-consuming recurrent
233 computations and constitute a well-controlled and difficult challenge for neural networks. It is well
234 known that humans can solve a number of visual tasks very quickly, presumably in a single
235 feedforward pass of neural activity (44). ffcNNs are good models of this kind of visual processing (17,
236 18, 45). However, neural activities are not determined by the feedforward sweep alone. Recurrent

activity is crucial for several reasons (6, 7, 46–49). First, information computed at a higher level can affect processing of local elements (for example, global configurations of flankers can affect processing of the local vernier target via feedback). Second, although feedforward networks can in principle implement any function (50), recurrent networks can implement these functions more efficiently, by recycling neural resources (48). Third, recurrent networks have the advantage of affording two distinct processing regimes (6): a fast feedforward pass able to quickly process information, and a time-consuming recurrent regime to perform more in depth global computations. CapsNets naturally include both a fast feedforward and a time-consuming recurrent regime. When a single routing by agreement iteration is used, CapsNets are rapid feedforward networks that can accomplish many tasks, such as vernier discrimination. With more routing iterations, a recurrent processing regime arises, and, with it, complex global shape effects emerge, such as computing and segmenting the cuboids in experiment 2. We showed how these two regimes in CapsNets explain our psychophysical results about temporal dynamics of (un)crowding by showing how recurrent processing kicks in when complex global processing is needed.

One limitation in our experiments is that we explicitly taught the CapsNets which configurations to group together by selecting which groups of shapes were present during training (e.g., only groups of identical shapes in experiment 1). Effectively, this gave the network adequate priors to produce uncrowding with the appropriate configurations (i.e., only identical, but not different flankers). Hence, our results show that, given adequate priors, CapsNets explain uncrowding. We have shown previously that ffCNNs do *not* produce uncrowding, *even* when they were similarly trained on groups of identical shapes and showed learning on the training data comparable to the CapsNets (15). This shows that merely training networks on groups of identical shapes is not sufficient to explain uncrowding. It is the recurrent segmentation in CapsNets that is crucial. Humans do not start from zero and therefore do not need to be trained in order to perform crowding tasks. The human brain is shaped through evolution and learning to group elements in a useful way to solve the tasks it faces. As mentioned, (un)crowding can be seen as a probe into this grouping strategy. Hence, we expect that training CapsNets on more naturalistic tasks such as ImageNet may lead to grouping strategies similar to humans and may therefore naturally equip the networks with priors that explain (un)crowding results. At the moment, however, CapsNets have not been trained on such difficult tasks because the routing by agreement algorithm is computationally expensive.

Recurrent networks are harder to train than feedforward systems, which explains the dominance of the latter during these early days of deep learning. However, despite this hurdle, recurrent networks

are emerging to address the limitations of ffCNNs as models of the visual system (7, 46, 48, 49, 51, 52). Our results suggest that one important role of recurrence is shape-level computations through grouping and segmentation. We had previously suggested another recurrent segmentation network, hard-wired to explain uncrowding (53). However, CapsNets, bringing together recurrent grouping and segmentation with the power of deep learning, are much more flexible and can be trained to solve any task. Linsley et al. (49) proposed another recurrent deep neural network for grouping and segmentation, and there are other possibilities too (54, 55). We do not suggest that CapsNets are the only implementation of grouping and segmentation.

In conclusion, our results provide mutually reinforcing modelling and psychophysical evidence that *time-consuming, recurrent* grouping and segmentation play a crucial role for global shape computations in humans. Recurrence kicks in when efficient grouping and segmentation of complex global shapes is required. We showed that CapsNets are a good model of this process. ffCNNs and other local feedforward models of vision, on the other hand, adopt a fundamentally different strategy for vision, which seems inadequate for human-like global shape computations.

283

284 Methods

The code to reproduce all our results will be available with the journal version of this contribution. All models were implemented in Python 3.6, using the high-level estimator API of Tensorflow 1.10.0. Computations were run on a GPU (NVIDIA GeForce GTX 1070). We used the same basic network architecture in all experiments (Figure 4a). We implemented early feature extraction by using three convolutional layers without padding, each followed by an ELU non-linearity. We used dropout (56) after the first and second convolutional layers. The outputs of the last convolution were reshaped into m primary capsule types outputting n -dimensional activation vectors. The number of output capsule types was equal to the number of different shapes used as input. The network was trained end-to-end through backpropagation. For training, we used an Adam optimizer with a batch size of 48 and a learning rate of 0.0004. To this learning rate, we applied cosine decays with warm restarts (57).

This choice of network architecture was motivated by the following rationale (Figure 4b). After training, ideally, primary capsules detect the individual shapes present in the input image, and output capsules group and segment these shapes through recurrent routing by agreement. The network can only group shapes together if it was taught during training that these shapes should form a group. To

match this rationale, we set the primary capsules' receptive field sizes to roughly the size of one shape, and we set the number of output capsules equal to the number of shape types. Inputs were grayscale images (Figure 4c&d). We added random Gaussian noise with mean $\mu = 0$ and standard deviation randomly drawn from a uniform distribution $\sigma \sim \mathcal{U}(0.00, 0.02)$. The contrast was varied either by first adding a random value between -0.1 and 0.1 to all pixel values and then multiplying them with a random value drawn from a uniform distribution $\mathcal{U}(0.6, 1.2)$, or vice versa. The pixel values were then clipped between 0 and 1.

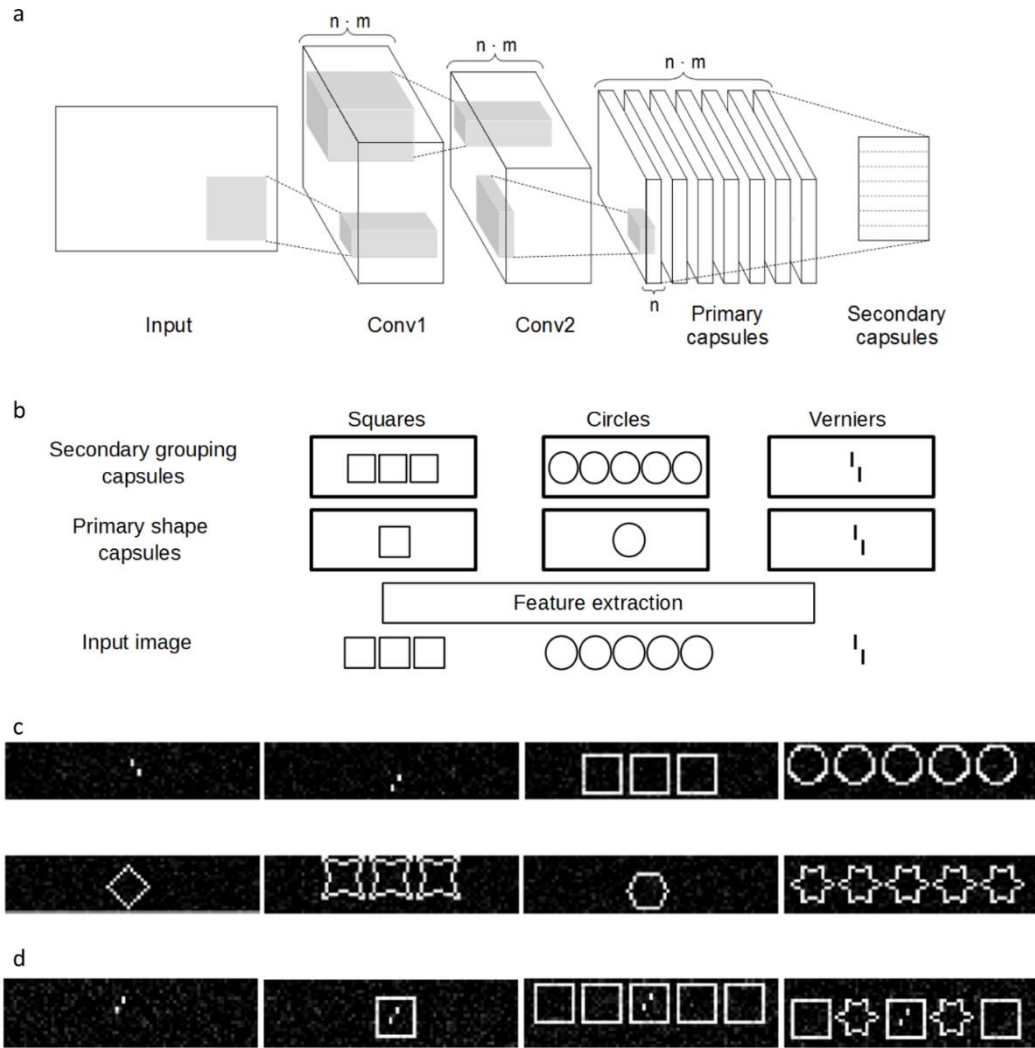


Figure 4: a. Network architecture: We used capsule networks with three convolutional layers whose last outputs was reshaped into the primary capsule layer with m primary capsule types and n primary capsule dimensions. In this example, the number of primary and output capsule types is seven to match the seven shape types we used in experiment 1 (see caption c), but the number depended on the experiment. The primary and output capsule layers communicate via routing-by-agreement. **b. Ideal representations:** After training, the primary capsules detect single shapes of different types at different locations. In this example, there are squares, circles and verniers. By routing the outputs of the primary capsules to the corresponding output capsules, the output capsules group these shapes in groups of one, three or five, based on the number of shapes detected by the primary capsules. If the left stimulus with three squares is presented, the primary

square capsules detect squares at three different locations. Through routing by agreement, the output squares capsule groups these three squares. If the middle stimulus with five circles is presented, the primary circle capsules detect circles at five different locations. Through routing by agreement, the output circles capsule represents a group of five circles after routing. Lastly, if a vernier is presented (right stimulus), it is detected by primary capsules and is represented in the vernier output capsule. **c. Training stimuli for experiment 1:** All shapes were shown randomly in groups of one, three or five, except verniers who were always presented alone. **d. Testing stimuli for experiment 1:** Example stimuli for the four test conditions: In the vernier-alone condition (*left*), we expected the network to perform well on the vernier discrimination task. In crowding conditions (*middle-left*), we expected a deterioration of the vernier discrimination as in classical crowding. In uncrowding conditions with many identical flankers (*middle-right*), we expected a recovery of the vernier discrimination. In no-uncrowding conditions with different flanker types (*right*), we expected crowding. After training, the network has learnt about groups of identical shapes and verniers, but has never encountered these (un)crowding stimuli.

Experiment 1:

Modelling

Human data for experiment 1 is based on (25). We trained CapsNets with the above architecture to solve a vernier offset discrimination task and classify groups of identical shapes. The training dataset included vernier stimuli and six different shape types (Figure 4c). Shapes were presented in groups of one, three or five shapes of the same type. The group was centered in the middle of the image, with a jitter of 2 pixels along the x-axis and 6 pixels along the y-axis.

The loss function included a term for shape type classification, a term for vernier offset discrimination, a term for the number of shapes in the image, and a term for reconstructing the input based on the network output (see equations 1-5). Each loss term was scaled so that none of the terms dominated the others. For the shape type classification loss, we implemented the same margin loss as in (2). This loss enables the detection of multiple objects in the same image. For the vernier offset loss, we used a small decoder to determine vernier offset directions based on the activity of the vernier output capsule. The decoder was composed of a single dense hidden layer followed by a ReLU-nonlinearity and a dense readout layer of two nodes corresponding to the labels left and right. The vernier offset loss was computed as the softmax cross entropy between the decoder output and the one-hot-encoded vernier offset labels. The loss term for the number of shapes in the image was implemented similarly, but the output layer comprised three nodes representing the labels one, three or five shape repetitions. For the reconstruction loss, we trained a decoder with two fully-connected hidden layers (h1: 512 units, h2: 1024 units) each followed by ELU nonlinearities to reconstruct the input image. The reconstruction loss was then calculated as the squared difference between the pixel values of the input image and the reconstructed image. The total loss is given by the following formulas:

$$L_{total} = \alpha_{shape\ type} L_{shape\ type} + \alpha_{vernier\ offset} L_{vernier\ offset} + \alpha_{shape\ repetitions} L_{shape\ repetitions} + \alpha_{reconstruction} L_{reconstruction} \quad (1)$$

$$L_{shape\ type} = \sum_k T_k \max(0, (m^+ - \|v_k\|)^2) + \lambda(1 - T_k) \max(0, (\|v_k\| - m^-)^2) \quad (2)$$

$$L_{vernier\ offset} = \text{Crossentropy}(\text{vernier labels}, \text{vernier decoder output}) \quad (3)$$

$$L_{shape\ repetitions} = \text{Crossentropy}(\text{shape repetitions labels}, \text{shape repetitions decoder output}) \quad (4)$$

$$L_{reconstruction} = \sum_{i,j} (\text{input}(i,j) - \text{reconstruction}(i,j))^2 \quad (5)$$

Where the α are real numbers scaling each loss term, $T_k = 1$ if shape class k is present, $\|v_k\|$ is the norm of output capsule k , and m^+ , m^- and λ are parameters of the margin loss with the same values as described in (2).

After training, we tested vernier discrimination performance on (un)crowding stimuli (figure 4d), and obtained input reconstructions. We trained 10 different networks and averaged their performance. Before this experiment, the network had never seen crowding nor uncrowding stimuli, but it knew about groups of shapes and about the vernier discrimination task. Therefore, the network could not trivially learn when to (un)crowd by overfitting on the training dataset. This situation is similar for humans: they know about shapes and verniers, but their visual system has never been trained on (un)crowding stimuli.

366

Experiment 2:

Psychophysical experiment:

Observers

For experiment 2, we collected human psychophysical data. Participants were paid students of the Ecole Polytechnique Fédérale de Lausanne (EPFL). All had normal or corrected-to-normal vision, with a visual acuity of 1.0 (corresponding to 20/20) or better in at least one eye, measured with the Freiburg Visual Acuity Test. Observers were told that they could quit the experiment at any time they wished. Five observers (two females) performed the experiment.

Apparatus and stimuli

376 Stimuli were presented on a HP-1332A XY-display equipped with a P11 phosphor and controlled by a
377 PC via a custom-made 16-bit DA interface. Background luminance of the screen was below 1 cd/m².
378 Luminance of stimuli was 80 cd/m². Luminance measurements were performed using a Minolta Lu-
379 minance meter LS-100. The experimental room was dimly illuminated (0.5 lx). Viewing distance was
380 75 cm.

381 We determined vernier offset discrimination thresholds for different flanker configurations. The ver-
382 nier target consisted of two lines that were randomly offset either to the left or right. Observers indi-
383 cated the offset direction. Stimulus consisted of two vertical 40' (arcmin) long lines separated by a
384 vertical gap of 4' and presented at an eccentricity of 5° to the right of a fixation cross (6' diameter).
385 Eccentricity refers to the center of the target location. Flanker configurations were centered on the
386 vernier stimulus and were symmetrical in the horizontal dimension. Observers were presented two
387 flanker configurations. In the lines configuration, the vernier was flanked by two vertical lines (84') at
388 40' from the vernier. In the cuboids configuration, perspective cuboids were presented to the left and
389 to the right of the vernier (width = 58', angle of oblique lines = 135°, length = 23.33'). Cuboids con-
390 tained the lines from the Lines condition as their centermost edge.

391 *Procedure*

392 Observers were instructed to fixate a fixation cross during the trial. After each response, the screen
393 remained blank for a maximum period of 3 s during which the observer was required to make a re-
394 sponse on vernier offset discrimination by pressing one of two push buttons. The screen was blank
395 for 500 ms between response and the next trial.

396 An adaptive staircase procedure (PEST; 58) was used to determine the vernier offset for which ob-
397 servers reached 75% correct responses. Thresholds were determined after fitting a cumulative Gauss-
398 ian to the data using probit and likelihood analyses. In order to avoid extremely large vernier offsets,
399 we restricted the PEST procedure to not exceed 33.3' i.e. twice the starting value of 16.66'. Each con-
400 dition was presented in separate blocks of 80 trials. All conditions were measured twice (i.e., 160
401 trials) and randomized individually for each observer. To compensate for possible learning effects, the
402 order of conditions was reversed after each condition had been measured once. Auditory feedback
403 was provided after incorrect or omitted responses.

404 **Modelling:**

405 To model the results of experiment 2, we trained our CapsNets to solve a vernier offset discrimination
406 task and classify verniers, cuboids, scrambled cuboids and lines. The training dataset included vernier
407 stimuli and one of three different shape types (lines, cuboids, scrambled cuboids). The scrambled

408 cuboids were included to make the task harder, and to prevent the network from classifying cuboids
409 simply based on the number of pixels in the image. The line stimuli were randomly presented in a
410 group of 2, 4, 6 or 8. Both, cuboids and shuffled cuboids were always presented in groups of two
411 facing one another. The distance between these shapes was varied randomly between one and six
412 pixels. The loss function was very similar to experiment 1, but without the loss term for shape repeti-
413 tions, since there were no repetitions (each term is the same as in eqs. 1-5):

$$414 \quad L_{total} = \alpha_{shape\ type} L_{shape\ type} + \alpha_{vernier\ offset} L_{vernier\ offset} + \alpha_{reconstruction} L_{reconstruction} \quad (6)$$

415 After training, we tested the network’s vernier discrimination performance on (un)crowding stimuli
416 (verniers surrounded by either lines, cuboids or scrambled cuboids), while varying the number of
417 recurrent routing by agreement iterations. We trained the same network 50 times and averaged per-
418 formance over these trained networks, excluding 21 networks for which vernier discrimination per-
419 formance with *both* line and cuboid flankers was at ceiling ($\geq 95\%$) or floor ($\leq 55\%$). This exclusion
420 criterion is used for cleaner results and does *not* impact the crucial result showing that uncrowding
421 occurs with increasing routing iterations only with cuboid, but not with line flankers. The effect still
422 occurs when all 50 networks are included in the analysis, but the fact that certain networks are at
423 floor or ceiling is misleading. Before this experiment, the network had never seen (un)crowding stim-
424 uli, but it knew about cuboids, scrambled cuboids and about the vernier discrimination task. There-
425 fore, the network could not trivially learn when to (un)crowd by overfitting on the training dataset.

426

427 Acknowledgements

428 Adrien Doerig was supported by the Swiss National Science Foundation grant n.176153 “Basics of
429 visual processing: from elements to figures”.

430

431 Bibliography

- 432 1. Doerig A, et al. (2019) Beyond Bouma’s window: How to explain global aspects of crowding?
433 *PLOS Computational Biology* 15(5):e1006580.
- 434 2. Sabour S, Frosst N, Hinton GE (2017) Dynamic routing between capsules. *Advances in Neural*
435 *Information Processing Systems*, pp 3856–3866.
- 436 3. DiCarlo JJ, Zoccolan D, Rust NC (2012) How Does the Brain Solve Visual Object Recognition?
437 *Neuron* 73(3):415–434.
- 438 4. Hubel DH, Wiesel TN (1962) Receptive fields, binocular interaction and functional architecture
439 in the cat’s visual cortex. *The Journal of physiology* 160(1):106–154.

- 440 5. Krizhevsky A, Sutskever I, Hinton GE (2012) Imagenet classification with deep convolutional
441 neural networks. *Advances in Neural Information Processing Systems*, pp 1097–1105.
- 442 6. Lamme VA, Roelfsema PR (2000) The distinct modes of vision offered by feedforward and
443 recurrent processing. *Trends in neurosciences* 23(11):571–579.
- 444 7. Kietzmann TC, et al. (2019) Recurrence required to capture the dynamic computations of the
445 human ventral visual stream. *arXiv preprint arXiv:190305946*.
- 446 8. Funke CM, et al. (2018) Comparing the ability of humans and DNNs to recognise closed contours
447 in cluttered images. *18th Annual Meeting of the Vision Sciences Society (VSS 2018)*, p 213.
- 448 9. Su J, Vargas DV, Sakurai K (2019) One pixel attack for fooling deep neural networks. *IEEE*
449 *Transactions on Evolutionary Computation*.
- 450 10. Szegedy C, et al. (2013) Intriguing properties of neural networks. *arXiv preprint arXiv:13126199*.
- 451 11. Geirhos R, et al. (2018) ImageNet-trained CNNs are biased towards texture; increasing shape
452 bias improves accuracy and robustness. *arXiv preprint arXiv:181112231*.
- 453 12. Baker N, Lu H, Erlikhman G, Kellman PJ (2018) Deep convolutional networks do not classify
454 based on global object shape. *PLoS computational biology* 14(12):e1006613.
- 455 13. Brendel W, Bethge M (2019) Approximating CNNs with Bag-of-local-Features models works
456 surprisingly well on ImageNet. *arXiv preprint arXiv:190400760*.
- 457 14. Kim T, Bair W, Pasupathy A (2019) Neural coding for shape and texture in macaque area V4.
458 *Journal of Neuroscience* 39(24):4760–4774.
- 459 15. Doerig A, Bornet A, Choung OH, Herzog MH (2019) Crowding Reveals Fundamental Differences
460 in Local vs. Global Processing in Humans and Machines. *bioRxiv*:744268.
- 461 16. Deng J, et al. (2009) Imagenet: A large-scale hierarchical image database. *2009 IEEE Conference*
462 *on Computer Vision and Pattern Recognition (Ieee)*, pp 248–255.
- 463 17. Khaligh-Razavi S-M, Kriegeskorte N (2014) Deep supervised, but not unsupervised, models may
464 explain IT cortical representation. *PLoS computational biology* 10(11):e1003915.
- 465 18. Yamins DL, et al. (2014) Performance-optimized hierarchical models predict neural responses in
466 higher visual cortex. *Proceedings of the National Academy of Sciences* 111(23):8619–8624.
- 467 19. RichardWebster B, Anthony S, Scheirer W (2018) Psyphy: A psychophysics driven evaluation
468 framework for visual recognition. *IEEE transactions on pattern analysis and machine*
469 *intelligence*.
- 470 20. Levi DM (2008) Crowding—An essential bottleneck for object recognition: A mini-review. *Vision*
471 *Research* 48(5):635–654.
- 472 21. Whitney D, Levi DM (2011) Visual crowding: a fundamental limit on conscious perception and
473 object recognition. *Trends in Cognitive Sciences* 15(4):160–168.

- 474 22. Bouma H (1973) Visual interference in the parafoveal recognition of initial and final letters of
475 words. *Vision Research* 13(4):767–782.
- 476 23. Pelli DG (2008) Crowding: a cortical constraint on object recognition. *Current Opinion in*
477 *Neurobiology* 18(4):445–451.
- 478 24. Manassi M, Whitney D (2018) Multi-level Crowding and the Paradox of Object Recognition in
479 Clutter. *Current Biology* 28(3):R127–R133.
- 480 25. Manassi M, Lonchampt S, Clarke A, Herzog MH (2016) What crowding can tell us about object
481 representations. *Journal of Vision* 16(3):35–35.
- 482 26. Westheimer G, Hauske G (1975) Temporal and spatial interference with vernier acuity. *Vision*
483 *research* 15(10):1137–1141.
- 484 27. Levi DM, Klein SA, Aitsebaomo AP (1985) Vernier acuity, crowding and cortical magnification.
485 *Vision research* 25(7):963–977.
- 486 28. Oberfeld D, Stahn P (2012) Sequential grouping modulates the effect of non-simultaneous
487 masking on auditory intensity resolution. *PloS one* 7(10):e48054.
- 488 29. Overvliet KE, Sayim B (2016) Perceptual grouping determines haptic contextual modulation.
489 *Vision Research* 126(Supplement C):52–58.
- 490 30. Saarela TP, Sayim B, Westheimer G, Herzog MH (2009) Global stimulus configuration modulates
491 crowding. *Journal of Vision* 9(2):5–5.
- 492 31. Herzog MH, Fahle M (2002) Effects of grouping in contextual modulation. *Nature*
493 415(6870):433.
- 494 32. Sayim B, Westheimer G, Herzog MH (2010) Gestalt factors modulate basic spatial vision.
495 *Psychological Science* 21(5):641–644.
- 496 33. Saarela TP, Westheimer G, Herzog MH (2010) The effect of spacing regularity on visual crowding.
497 *Journal of Vision* 10(10):17–17.
- 498 34. Manassi M, Sayim B, Herzog MH (2012) Grouping, pooling, and when bigger is better in visual
499 crowding. *Journal of Vision* 12(10):13–13.
- 500 35. Pachai MV, Doerig AC, Herzog MH (2016) How best to unify crowding? *Current Biology*
501 26(9):R352–R353.
- 502 36. Manassi M, Sayim B, Herzog MH (2013) When crowding of crowding leads to uncrowding.
503 *Journal of Vision* 13(13):10–10.
- 504 37. Wallace JM, Chiu MK, Nandy AS, Tjan BS (2013) Crowding during restricted and free viewing.
505 *Vision Research* 84:50–59.
- 506 38. Tripathy SP, Cavanagh P, Bedell HE (2014) Large crowding zones in peripheral vision for briefly
507 presented stimuli. *Journal of Vision* 14(6):11–11.

- 508 39. Styles EA, Allport DA (1986) Perceptual integration of identity, location and colour. *Psychological*
509 *Research* 48(4):189–200.
- 510 40. Li Z (1999) Visual segmentation by contextual influences via intra-cortical interactions in the
511 primary visual cortex. *Network: computation in neural systems* 10(2):187–212.
- 512 41. Parkes L, Lund J, Angelucci A, Solomon JA, Morgan M (2001) Compulsory averaging of crowded
513 orientation signals in human vision. *Nature neuroscience* 4(7):739.
- 514 42. Pelli DG, Palomares M, Majaj NJ (2004) Crowding is unlike ordinary masking: Distinguishing
515 feature integration from detection. *Journal of Vision* 4(12):12–12.
- 516 43. Rosenholtz R, Yu D, Keshvari S (2019) Challenges to pooling models of crowding: Implications for
517 visual mechanisms. *Journal of vision* 19(7).
- 518 44. Thorpe S, Fize D, Marlot C (1996) Speed of processing in the human visual system. *nature*
519 381(6582):520.
- 520 45. Kietzmann TC, McClure P, Kriegeskorte N (2018) Deep neural networks in computational
521 neuroscience. *bioRxiv*:133504.
- 522 46. Kim J, Linsley D, Thakkar K, Serre T (2019) Disentangling neural mechanisms for perceptual
523 grouping. *arXiv preprint arXiv:190601558*.
- 524 47. Tang H, et al. (2018) Recurrent computations for visual pattern completion. *Proceedings of the*
525 *National Academy of Sciences* 115(35):8835–8840.
- 526 48. Spoerer CJ, Kietzmann TC, Kriegeskorte N (2019) Recurrent networks can recycle neural
527 resources to flexibly trade speed for accuracy in visual recognition. *bioRxiv*:677237.
- 528 49. Linsley D, Kim J, Serre T (2018) Sample-efficient image segmentation through recurrence.
529 *arXiv:181111356 [cs]*. Available at: <http://arxiv.org/abs/1811.11356> [Accessed June 27, 2019].
- 530 50. Hornik K, Stinchcombe M, White H (1989) Multilayer feedforward networks are universal
531 approximators. *Neural networks* 2(5):359–366.
- 532 51. Spoerer CJ, McClure P, Kriegeskorte N (2017) Recurrent convolutional neural networks: a better
533 model of biological object recognition. *Frontiers in psychology* 8:1551.
- 534 52. Kar K, Kubilius J, Schmidt K, Issa EB, DiCarlo JJ (2019) Evidence that recurrent circuits are critical
535 to the ventral stream’s execution of core object recognition behavior. *Nature neuroscience*
536 22(6):974.
- 537 53. Francis G, Manassi M, Herzog MH (2017) Neural dynamics of grouping and segmentation
538 explain properties of visual crowding. *Psychological review* 124(4):483.
- 539 54. Ronneberger O, Fischer P, Brox T (2015) U-net: Convolutional networks for biomedical image
540 segmentation. *International Conference on Medical Image Computing and Computer-Assisted*
541 *Intervention* (Springer), pp 234–241.
- 542 55. Girshick R, Radosavovic I, Gkioxari G, Dollár P, He K (2018) *Detectron*.

543 56. Srivastava N, Hinton G, Krizhevsky A, Sutskever I, Salakhutdinov R (2014) Dropout: a simple way
544 to prevent neural networks from overfitting. *The Journal of Machine Learning Research*
545 15(1):1929–1958.

546 57. Loshchilov I, Hutter F (2016) Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint*
547 *arXiv:160803983*.

548 58. Taylor M, Creelman CD (1967) PEST: Efficient estimates on probability functions. *The Journal of*
549 *the Acoustical Society of America* 41(4A):782–787.

550

551

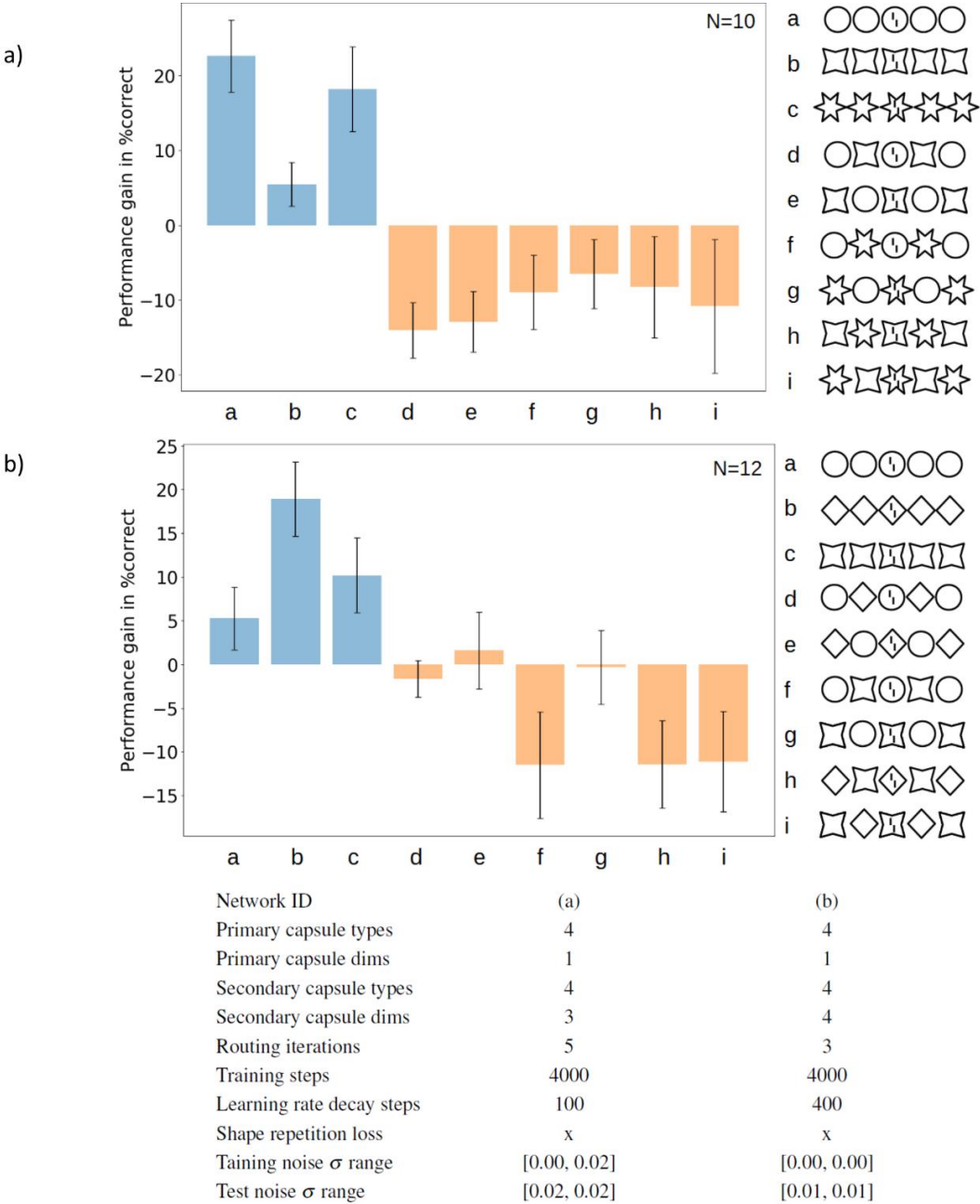
Supplementary Material

Experiment 1

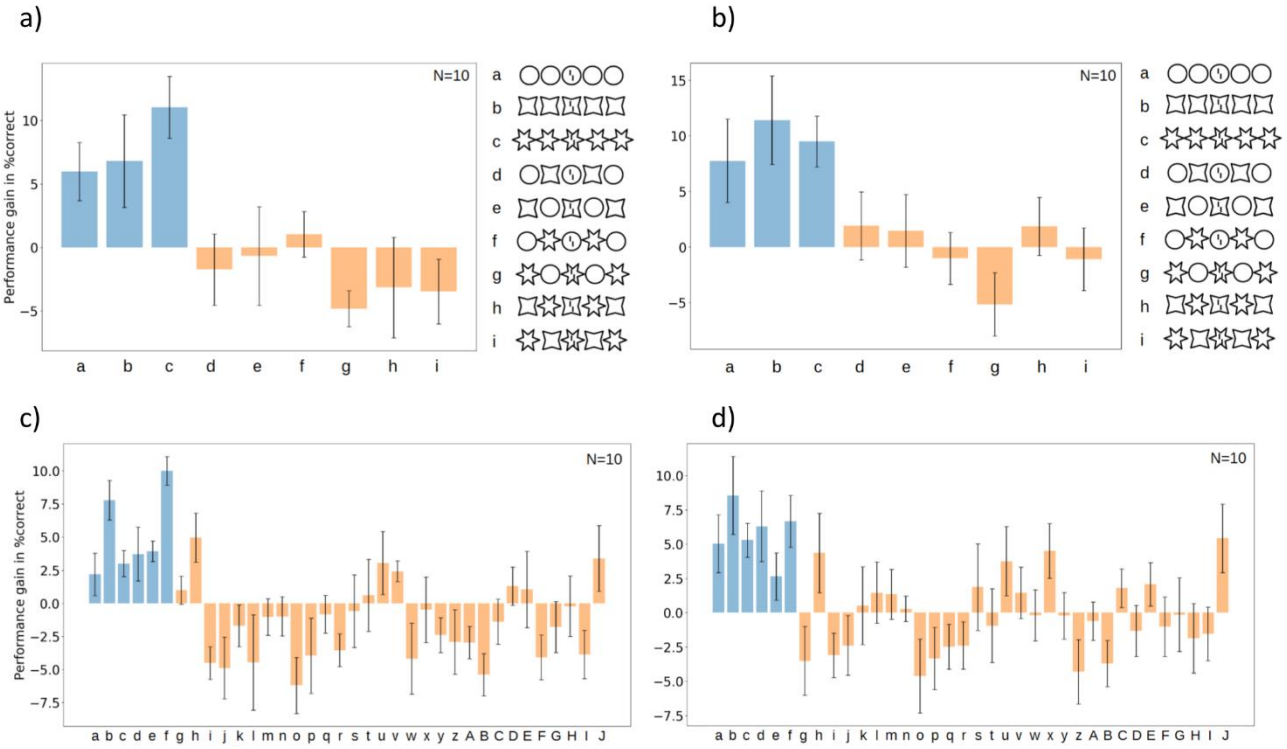
Results are robust against stimuli and hyperparameters changes

To avoid cherrypicking our hyperparameters, we ran several networks with different hyperparameter sets, and show that our results are robust with respect to these changes.

The results of experiment 1 remain qualitatively similar for different image sizes and network hyperparameters. Below is a selection of results using different sets of hyperparameters. In all these cases, both crowding and uncrowding occur, similarly to the results shown in Figure 2.

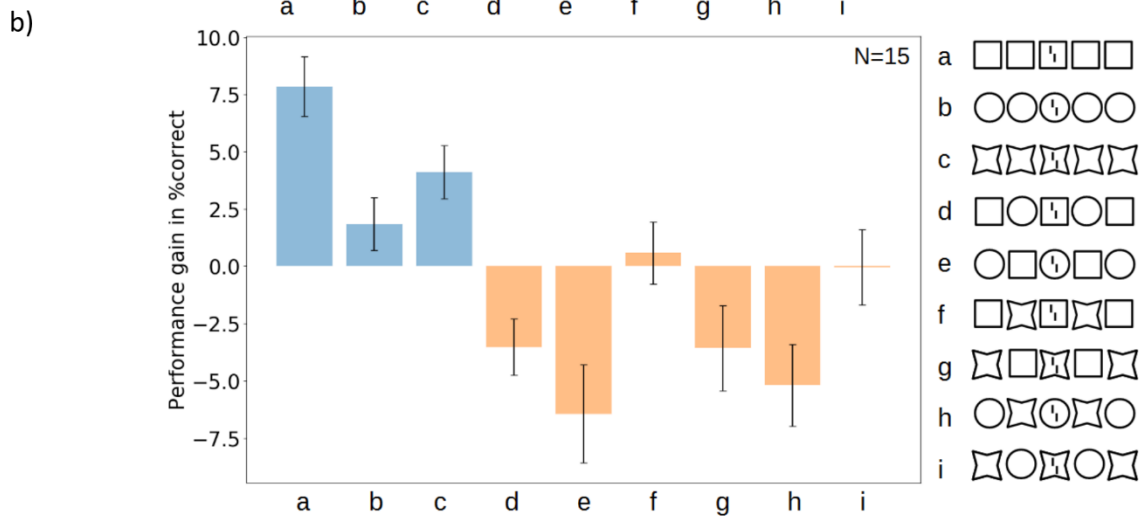
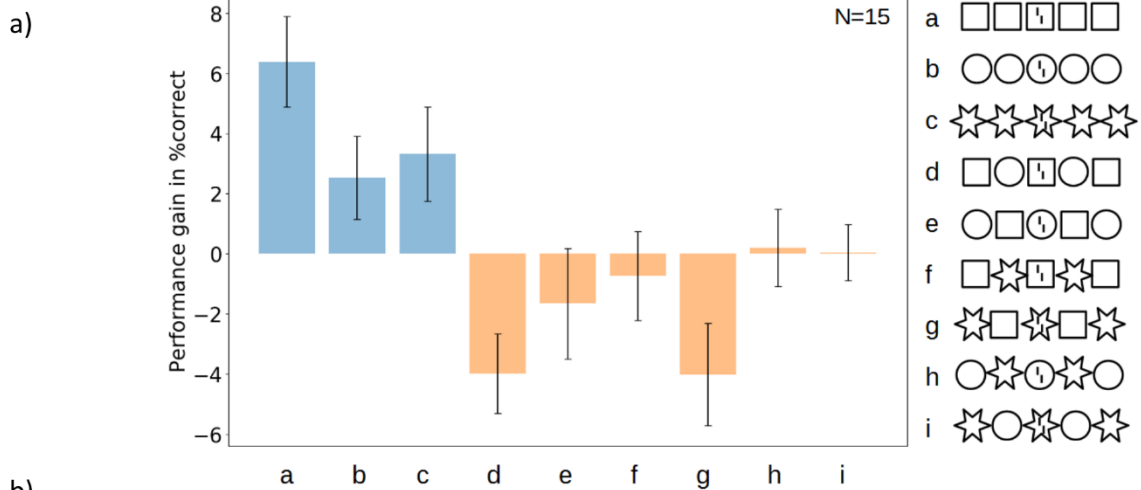


Supplementary Figure 1: Results for 16x72 pixel images. Both crowding and uncrowding occur similarly to the results in figure 2. Plotting conventions are the same as in figure 2. Main hyperparameters are summarized at the bottom. With these small images, we often encountered ceiling effects. We trained 20 networks and dropped those that were at ceiling (i.e., we dropped networks that were at 100% performance for all conditions).



Network ID	(a)	(b)	(c)	(d)
Primary capsule types	4	4	7	7
Primary capsule dims	1	1	2	2
Secondary capsule types	4	4	7	7
Secondary capsule dims	4	4	8	10
Routing iterations	3	5	3	3
Training steps	8000	6000	2500	5000
Shape repetition loss	x	x	x	x
Location loss			x	x
Reconstruction loss			x	x
Gaussian training noise	[0.00, 0.05]	[0.00, 0.00]	[0.02, 0.04]	[0.02, 0.04]
Gaussian test noise	[0.05, 0.05]	[0.05, 0.05]	[0.04, 0.06]	[0.04, 0.06]

Supplementary Figure 2: 20x72 pixel images. Both crowding and uncrowding occur similarly to the results in figure 2. Plotting conventions are the same as in figure 2. Main hyperparameters are summarized at the bottom. Stimuli not shown for panels b&c, for clarity.



Network ID	(a)	(b)
Primary capsule types	20	20
Primary capsule dims	1	1
Secondary capsule types	4	4
Secondary capsule dims	12	12
Routing iterations	4	5
Training steps	3000	3000
First decay steps	500	500
Shape repetition loss	x	x
Gaussian training noise	[0.00, 0.00]	[0.00, 0.02]
Gaussian test noise	[0.01, 0.01]	[0.01, 0.01]

571

572 **Supplementary Figure 3: 30x72 pixel images.** Both crowding and uncrowding occur similarly to the results in figure 2.

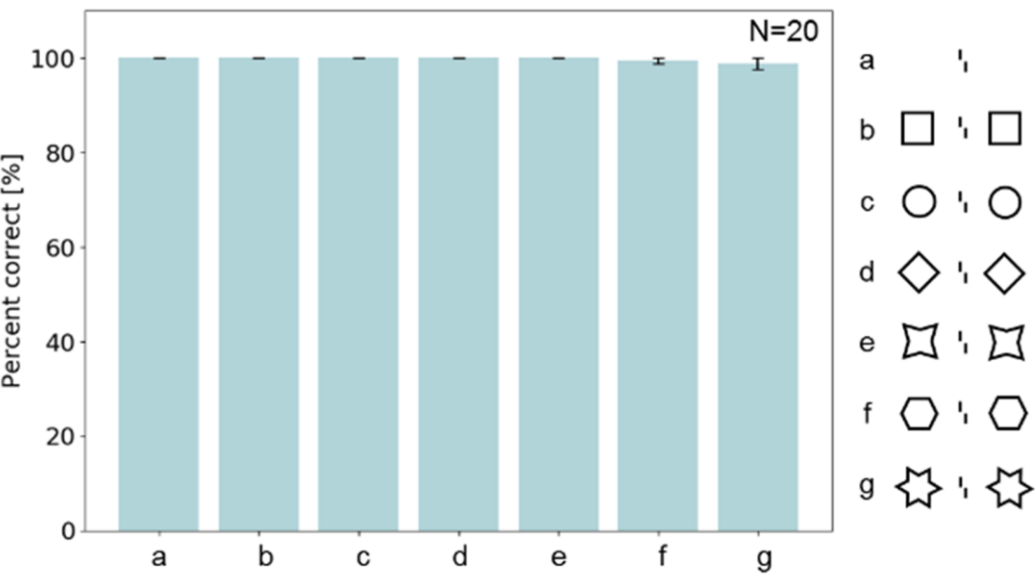
573 Plotting conventions are the same as in figure 2. Main hyperparameters are summarized at the bottom.

574

575 *Performance deterioration is due to crowding*

576 As a control to check that performance dropped because of crowding and not merely because of the
 577 simultaneous presentation of a vernier target and another shape, we measured performance when
 578 the vernier was presented outside, rather than inside, flanking shapes. Performance does not drop in

579 this case, compared to when the vernier is presented alone. This suggests that performance drops
 580 because of crowding in the networks.



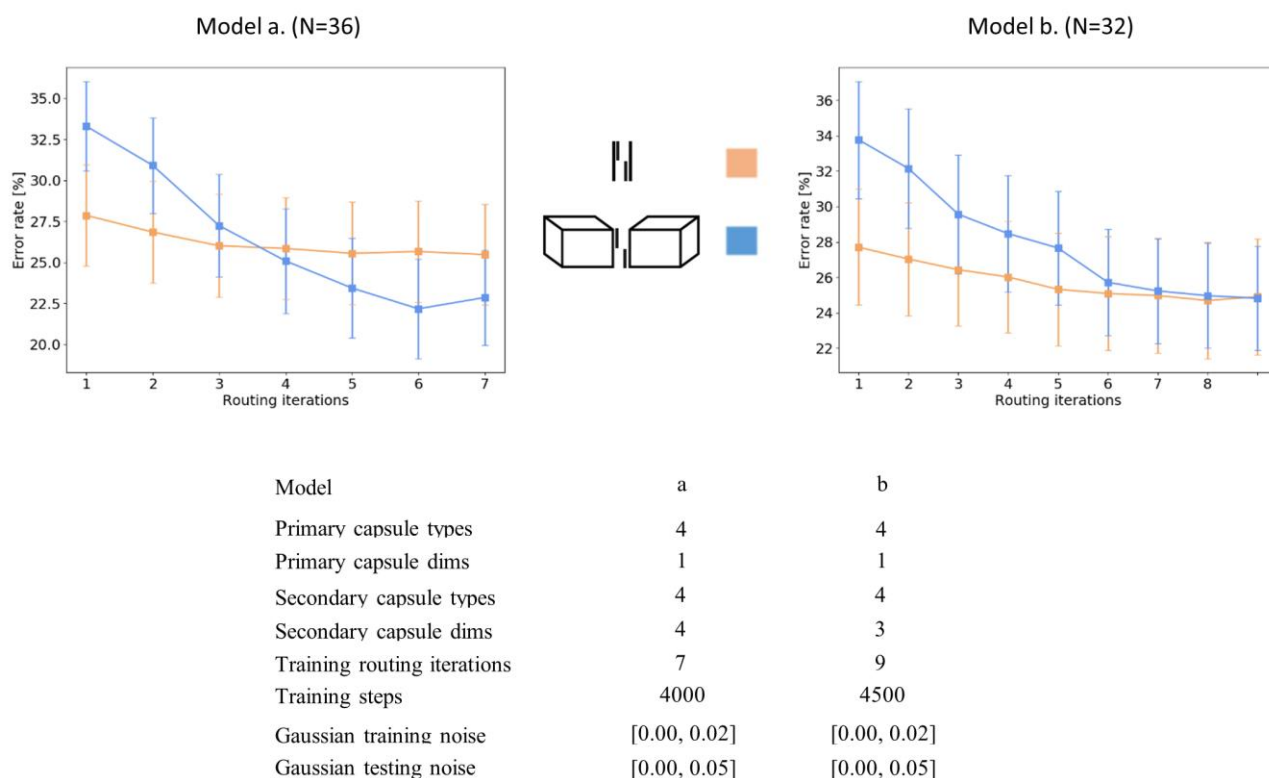
581
 582 **Supplementary Figure 4: Performance deterioration is due to crowding.** The x-axis shows different conditions shown on
 583 the right, the y-axis shows vernier offset discrimination percent correct. Vernier accuracy does not decrease when the
 584 vernier is presented outside flanking shapes compared to the vernier alone condition.

585
 586 **Experiment 2**

587 *Results are robust against stimuli and hyperparameters changes*

588 To avoid cherrypicking our hyperparameters, we ran several networks with different hyperparameter
 589 sets, and show that our results are robust with respect to these changes.

590 The results of experiment 2 remain qualitatively similar for different network hyperparameters. Below
 591 is a selection of results using different sets of hyperparameters. In both these cases, performance on
 592 the cuboids condition, but not the lines condition, drastically improves with the number of recurrent
 593 routing by agreement iterations (network a: lines: $p = 0.041$ vs. cuboids $p = .0.0005$, network b: lines:
 594 0.11 vs. cuboids $p=0.006$). In network a, the lines show a marginally significant improvement, but the
 595 p-value is 100 times smaller than for the cuboids.



596

597
598
599
600
601
602
603
604
605
606

Supplementary Figure 5: Experiment 2 results are reproduced with different network hyperparameters. The x-axis shows different numbers of routing iterations during testing and the y-axis shows the corresponding error rates (i.e., lower values indicate better performance). Error bars indicate standard deviation across N trained networks (see Methods). Performance increases drastically with recurrent routing iterations only for the cuboids condition, and not for the lines condition. A difference with the results shown in figure 3 is that performance with cuboids flankers is worse than performance with line flankers at early iterations. This may be explained by the far greater amount of pixels in cuboids than lines, increasing the interference between the cuboids and the vernier until the cuboids are segmented away. As the results exhibited in Figure 3 show, this effect can be mitigated through adequate hyperparameter choice. However, in this experiment, we focused on demonstrating that only the cuboids benefit from additional routing iterations, and this result is very stable across hyperparameter changes.