

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32

Title

- Brain-like border ownership signals support prediction of natural videos

Authors

Zeyuan Ye¹, Ralf Wessel¹, Tom P. Franken^{2*}

Affiliations

¹Department of Physics, Washington University in St. Louis, St. Louis, Missouri, USA

²Department of Neuroscience, Washington University School of Medicine, St. Louis, Missouri, USA

*Correspondence and requests for materials should be addressed to T.P.F. (email: ftom@wustl.edu)

Abstract

To make sense of visual scenes, the brain must segment foreground from background. This is thought to be facilitated by neurons in the primate visual system that encode border ownership (BOS), i.e. whether a local border is part of an object on one or the other side of the border. It is unclear how these signals emerge in neural networks without a teaching signal of what is foreground and background. In this study, we investigated whether BOS signals exist in PredNet, a self-supervised artificial neural network trained to predict the next image frame of natural video sequences. We found that a significant number of units in PredNet are selective for BOS. Moreover these units share several other properties with the BOS neurons in the brain, including robustness to scene variations that constitute common object transformations in natural videos, and hysteresis of BOS signals. Finally, we performed ablation experiments and found that BOS units contribute more to prediction than non-BOS units for videos with moving objects. Our findings indicate that BOS units are especially useful to predict future input in natural videos, even when networks are not required to segment foreground from background. This suggests that BOS neurons in the brain might be the result of evolutionary or developmental pressure to predict future input in natural, complex dynamic visual environments.

33

34 MAIN TEXT

35

36 Introduction

37 To understand the world around us, we parse incoming visual information into an organized collection
38 of objects. In primate animals, this capability is thought to be facilitated by neurons in the early areas in
39 visual cortex that encode border ownership (BOS)¹⁻⁴. These neurons fire more to an identical border in
40 their classical receptive field (cRF) depending on which side owns the border, even though the
41 contextual information that defines the side of foreground occurs far outside of the cRF (Figure 1A).
42 This selectivity extends to natural images^{5,6} and the preferred side of ownership corresponds to the side
43 that is near when varying depth⁷. Psychophysics and imaging studies support that BOS neurons also
44 exist in the human brain⁸⁻¹¹. It is unknown under which conditions BOS signals emerge in neural
45 networks. Artificial neural networks (ANNs) are a great tool to study such ‘why’ questions of how the
46 brain works, because they enable to test whether a particular neural phenomenon results from
47 optimization for a specific task¹².

48 It seems intuitive to hypothesize that BOS signals emerge in ANNs when they are explicitly trained on
49 scene segmentation, given that this is assumed to be the primary role of such signals in the brain¹³. A
50 recent study indeed found that units selective for BOS occur in a supervised ANN trained to segment
51 handwritten digits (a processed MNIST dataset¹⁴) from background¹⁵. However, such supervised
52 learning has been criticized as biologically highly implausible because it requires a large number of
53 explicitly segmented labels which is unrealistic in brain development^{16,17}. Another study found that BOS
54 signals can arise in an unsupervised ANN trained to develop translation invariance for an object
55 presented in isolation, but this mechanism failed for scenes with more than one object¹⁸, as opposed to
56 BOS signals in the brain¹⁹. Furthermore, these ANNs can only process simple artificial datasets, unlike
57 neural networks in modern deep-learning frameworks or the human brain, which are high performing on
58 realistic natural visual inputs²⁰⁻²². It thus remains poorly understood when BOS signals emerge in neural
59 networks.

60 Certain properties of BOS neurons in the brain suggest that BOS signals may be important under
61 dynamic conditions. BOS signals are known to persist for hundreds of milliseconds when the contextual
62 information that defines the side of ownership disappears, as long as the border in the cRF, which has
63 then become ambiguous for BOS, remains^{23,24}. Furthermore, these persistent BOS signals can be
64 transferred to other neurons if the ambiguous border lands in their cRF after an eye movement²⁵. This
65 hysteresis may make it easier to make sense of dynamic visual input by providing spatiotemporal
66 contiguity.

67 These observations motivated us to study whether BOS signals emerge in an artificial neural network
68 trained to predict future visual input for natural videos. We studied PredNet, a deep neural network with
69 an architecture inspired by predictive coding^{17,26,27}. PredNet was trained on a dataset of natural videos
70 captured by car-mounted cameras (KITTI²⁸) to predict the next video frame. Our *in-silico* experiments
71 demonstrate that a significant fraction of units in PredNet exhibit BOS signals. Moreover, these BOS
72 units share several properties with BOS neurons in the brain. Finally, ablating PredNet’s BOS units
73 increased prediction error more than ablating the same number of non-BOS units. BOS units thus
74 contribute to prediction of natural visual input even if there is no need to segment foreground from
75 background. This suggests that the need to predict future input in natural videos may drive the
76 development of BOS neurons. These scene segmentation signals, typically considered an example of a

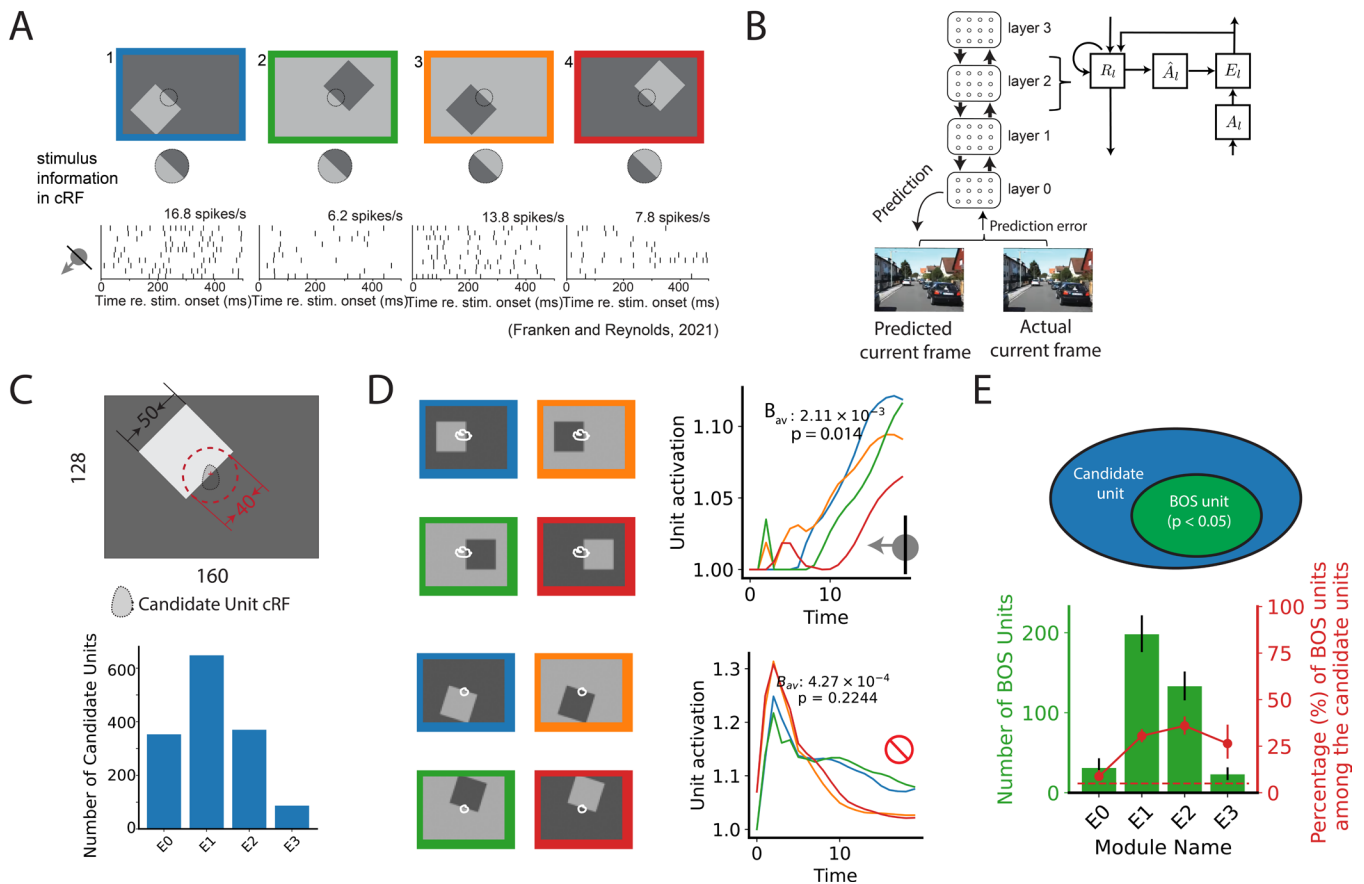
77 ventral ‘what’ stream operation, may thus be more involved in processing dynamic aspects of visual
78 input than is typically assumed.
79

80 **Results**

81 **BOS signals emerge in an artificial network trained to predict the next frame of natural videos**

82 To study the role of BOS signals in the processing of complex dynamic input we employed PredNet, a
 83 hierarchical ANN introduced by Lotter et al.²⁶ (Figure 1B). PredNet comprises four layers, with four
 84 modules per layer: the representation (R_l), the predicted output (\hat{A}_l), the prediction target (A_l), and the
 85 prediction error (E_l) modules. At each time step (see Methods), signals propagate from the top layer to
 86 the zeroth layer, resulting in a prediction for the next video frame in \hat{A}_0 . This prediction is then
 87 compared to the actual next frame provided in A_0 . The prediction error signal subsequently propagates
 88 from the zeroth layer to the top layer. The network was trained to minimize the prediction error of
 89 videos in the KITTI dataset²⁸, which were captured by car-mounted cameras in various urban and rural
 90 settings in Germany.

91
 92



93 **Figure 1. Border ownership (BOS) signals emerge in PredNet (A)** An example unit in the primate visual
 94 cortex that is selective for BOS. The unit has different responses depending on the BOS, even though the image
 95 pixels in its classical receptive field (cRF) are identical for panel 1 and 2, and for panel 3 and 4. The preferred
 96 side of ownership is the same for borders in the cRF with a different contrast polarity (the unit fires more to scene 1
 97 than to scene 2, and more to scene 3 than to scene 4). Arrow on the bottom left indicates the side of BOS that this
 98 unit prefers. Figure adapted from Franken and Reynolds³. **(B)** PredNet is an artificial neural network designed for
 99 video prediction. At each time step, the model operates by updating unit activities sequentially from the top layer
 100 (layer 3) to the bottom layer (layer 0), generating a prediction of the current video frame. The prediction error is
 101

102 then fed forward to layer 3. Each layer contains four modules (\hat{A}_l, A_l, E_l , and R_l where $l = 0, 1, 2, 3$ indicates
103 layer index, see Methods) (C) Candidate units in PredNet are defined as units whose cRF overlaps with the
104 central border but not with any of the square's corners (see Methods). Bottom: the number of candidate units in E
105 modules across different layers. See SI Figure 2 for R module data. (D) Responses of two example units (module
106 E_2), with white contours indicating the cRF. B_{av} measures, for each unit, the selectivity for BOS across different
107 square orientations (see Methods). Colored lines indicate the response to the different stimulus conditions (colors
108 indicate for each response function to which of the stimulus panels on the left it corresponds) for one orientation.
109 p value (two-tailed) was computed by comparing B_{av} to that obtained by shuffling the BOS labels (permutation
110 test, see Methods). BOS units are defined as those units for which this value is smaller than 0.05. (E) Number and
111 percentage of BOS units in E-module in different layers. Error bars indicate 95% confidence interval. Horizontal
112 dashed line indicates the chance level for the percentage of BOS units (5%).
113

114 We tested whether the BOS units exist in the PredNet by doing an *in-silico* experiment that is analogous
115 to the neurophysiological studies on BOS^{1,3,29} (Figure 1A). We measured the cRF of each unit using
116 sparse noise stimuli (SI Figure 1). First we identified candidate units for BOS tuning. For a unit to be a
117 candidate unit, the cRF needed to include the center of the square's central border (i.e. the border
118 positioned at the scene center) but must exclude any other border of the square (Fig. 1C; Methods). This
119 criterion is similar to that used in neurophysiological studies on BOS^{1,3}. We found tens to hundreds of
120 candidate units in different PredNet modules (Fig. 1C and SI Figure 2). We then analyzed the response
121 from candidate units to the standard full square stimuli. Figure 1D shows responses from two example
122 candidate units. The top unit exhibited a larger response when the square was positioned on the left side
123 of the central vertical border compared to the right side, irrespective of the contrast polarity across the
124 border (i.e. blue vs. green, and orange vs. red). This unit thus prefers that the border in its cRF is owned
125 by a square on one side, similar to BOS neurons in the primate visual cortex^{1,30}. In contrast, the bottom
126 unit did not exhibit a clear difference: its response was very similar for stimuli with opposite border
127 ownership but identical contrast polarity of the central border. To quantify BOS tuning for each unit, we
128 first computed the unit's difference in response between stimuli of opposite border ownership across
129 different contrast polarities, and divided it by the sum of the responses, resulting in the BOS index
130 (BOI). The BOS index was then averaged across all square orientations, resulting in B_{av} (see Methods).
131 The statistical significance of B_{av} was determined by comparing it to a null distribution obtained by
132 shuffling the stimulus labels. A candidate unit with a p-value smaller than 0.05 was defined as a BOS
133 unit; otherwise, it was defined as a non-BOS unit. The top unit in Figure 1D has a statistically significant
134 B_{av} and is therefore a BOS unit, while the bottom unit is a non-BOS unit.

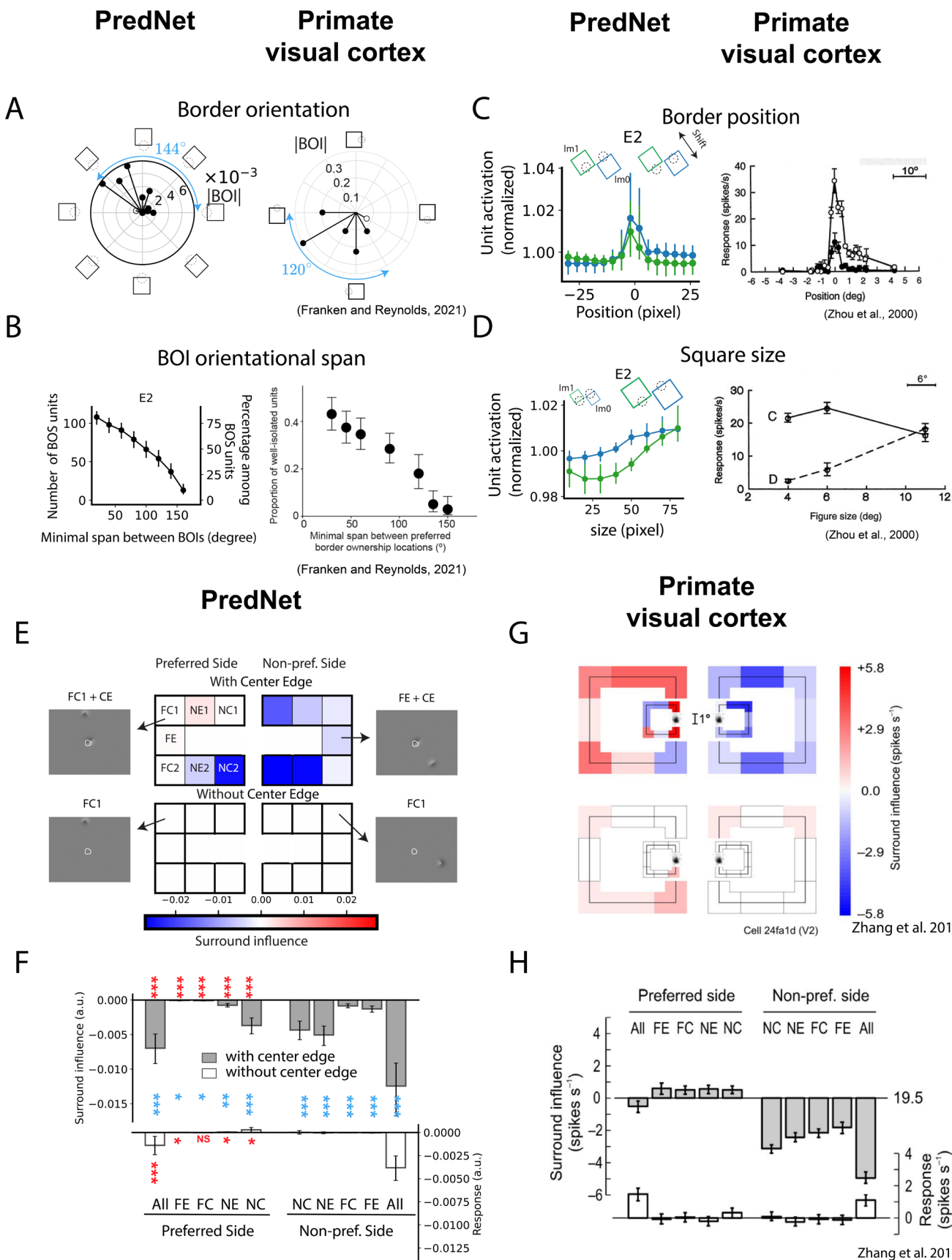
135 We conducted a population analysis of B_{av} across all candidate units. We find that 20-40% of candidate
136 units in E_1, E_2 and E_3 have significant B_{av} values (Fig. 1E). This is larger than in module E_0 (8.7 %,
137 95% confidence interval [6.2%, 12%]). BOS units were also found in modules R_1 and R_2 , but much less
138 in R_0 (R_3 had only a small number of candidate units; SI Fig. 2). The distribution of BOS units in
139 PredNet's hierarchy is reminiscent of the distribution of BOS neurons in the primate visual cortex,
140 which are less prevalent in areas closer to the sensory input (V1) than in downstream areas (V2 and
141 V4)^{1,3,4}.

142 **PredNet's BOS signals are robust to scene variations common in natural object** 143 **transformations, like BOS neurons in the brain**

144 We explored the robustness of BOS signals to the same scene variations that have been used in
145 neurophysiology studies on BOS neurons: square orientation, position, and size^{1,3}. Figure 2A (left)

146 shows the BOI for different square orientations for an example BOS unit. Vector length indicates the
147 absolute value of the BOI, and the angle of each vector indicates the preferred side of BOS for each
148 orientation (cf. the symbols around the plot). The square orientations with a large BOI form a contiguous
149 region in visual space, which is similar to BOS neurons in the primate visual cortex (e.g. Fig. 2A,
150 right)³. Filled symbols in Fig. 2A indicate orientations for which BOI is significant (permutation test, see
151 Methods). The angular span between object locations at the preferred side of BOS for different border
152 orientations (only orientations for which BOIs is significant are considered) is referred to as the BOS
153 span. For example, the span for the neuron shown in Figure 2A (left panel) is 144°. A substantial
154 number of BOS units in PredNet have a large span, extending to ~150°, similar to BOS neurons in the
155 brain (Fig. 2B; SI Fig. 4A).

156 Next, we examined BOS tuning for different square positions and sizes. We set the orientation at that for
157 which |BOI| was maximal, and then varied position (the position of the center of the central border
158 varied along a line orthogonal to the border's orientation; Fig. 2C). We also varied square size for the
159 central position (Fig. 2D). We find that the response difference between scenes with opposite BOS was
160 consistent for different positions or sizes in the population of BOS units in PredNet (Fig. 2C,D, left
161 panels), just like for BOS neurons in the brain (Fig. 2C,D, right panels)¹. We quantified this consistency
162 by averaging BOI across different conditions (i.e., size or position, see Methods). For all modules with a
163 substantial number of BOS units (over 15 units), these averaged BOI values are statistically significantly
164 positive (i.e. consistent with the tuning in the baseline condition; SI Figure 4B, C, bootstrapping test, see
165 Methods). Taken together, we find that border ownership signals in PredNet are robust to differences in
166 square orientation, size and position, i.e. remarkably similar to BOS neurons in the primate visual
167 cortex^{1,3}.



169 **Figure 2. BOS units in PredNet share several properties with BOS neurons in the brain. (A)** Border
170 ownership index (BOI) at different orientations. Vector magnitude represents the absolute value of BOI, i.e. the
171 difference in unit response to scenes with squares that share a border with a given orientation, but for which the
172 square is positioned on opposite sides of that border (thus a pair of stimulus cartoons on opposite sides of the
173 polar plot), divided by their sum (see Methods). Vector angle is such that the vector points towards the stimulus
174 cartoon with the preferred side of BOS for that border orientation. Left: example of a BOS unit in PredNet. Right:
175 BOS neuron recorded in macaque area V4 (reproduced from Franken and Reynolds 2021³). Filled symbols in
176 both panels indicate for which orientations the BOI was significantly different from 0. Blue text indicates the span
177 of a unit, which is the angle between the preferred object locations for orientations with statistically significant
178 BOI (Methods). **(B)** The y-axis illustrates the number/percentage of BOS units whose spans equal or exceed the
179 span indicated by the x values. Error bars indicate 95% confidence intervals. Left: BOS units in module E_2 in
180 PredNet (see SI Fig. 4A for other modules). Right: population data from BOS neurons in macaque area V4
181 (reproduced from Franken and Reynolds 2021³). **(C)** Left: for each BOS unit, squares were generated with
182 different positions as indicated in the cartoon. The blue and green traces represent normalized population
183 responses (see Methods) to opposite BOS (blue corresponds to the preferred BOS derived from responses to the
184 standard square set). Dots and error bars show the median, first, and third quantiles across the population of BOS
185 units. Right: responses from a BOS neuron in macaque V2 for different square positions. The two traces indicate
186 opposite BOS. Dots and error bars represent mean firing rates and SEMs across trials (reproduced with
187 permission from Zhou et al. 2000¹. Copyright 2000 Society for Neuroscience). **(D)** Identical to (C) but square size
188 was varied instead of square position. Right panel reproduced with permission from Zhou et al. 2000¹. Copyright
189 2000 Society for Neuroscience. **(E)** Response of an example BOS unit in PredNet (module E_2) to square
190 fragments. Top half shows responses to a square fragment in the surround paired with the border in the cRF.
191 Bottom half shows responses to square fragments in the surround without the border in the cRF. Gray panels
192 show example scenes (white outline: cRF). Colors of the central panels indicate the surround influence. The
193 surround influence is the unit's response to a scene with a square fragment in the surround at the position
194 indicated by the letter codes (also symbolized by the panel's position), subtracted by that to a scene without the
195 square fragment. Letter codes: NC: near corner; NE: near edge; FC: far corner; FE: far edge; numbers indicate
196 different positions of the fragment, e.g. NC1 and NC2 refer to each of the near corners on opposite ends of the
197 central border. **(F)** Means and 95% confidence intervals (i.e. 1.96 times SEM) of surround influence across all
198 BOS units in module E_2 ($n=71$; ; see SI Fig. 5 for other modules). NC is the average of NC1 and NC2, and the
199 same was done for NE and FC. 'All' represents the surround influence when all square segments were shown
200 (top), or all square segments except the center edge were show (bottom). Red text indicates whether the surround
201 influence for a particular condition is significantly larger on the preferred side than on the non-preferred side.
202 Blue text indicates whether the absolute value of surround influence of with-CE scenes is significantly larger than
203 without-CE scenes. Wilcoxon signed-rank test. ***: $p < 0.001$; **: $p < 0.01$; *: $p < 0.05$; NS: no significance.
204 Outlier units (see Methods) were removed to compute mean and SEM but included in the statistical tests. **(G)**
205 Same as (E) for a BOS neuron in the macaque visual cortex (reproduced with permission from Zhang et al.
206 2010³¹). Two different square sizes were evaluated, for which surround influence is plotted separately as the
207 smaller and larger panels. **(H)** Similar panel as (F), for BOS neurons in the macaque visual cortex, with
208 permission from Zhang et al. 2010³¹.

209 Surround influence for PredNet's BOS units is similar to BOS neurons in the brain

210 Neurophysiological experiments found that isolated object fragments in the surround modulate the
211 activity of BOS neurons in a way that is consistent with BOS tuning: fragments on the non-preferred
212 side of BOS suppress the response significantly more than fragments on the preferred side, which often
213 have an enhancing effect³¹. These modulatory effects were only significant in the presence of a border in
214 the cRF. We analyzed how fragments in the surround modulated the activity of BOS units in PredNet.

215 Similar to Zhang et al. 2010³¹, we divided a square object into 8 fragments: one Center Edge (CE)
216 located at the image center, and 7 contextual fragments (two Near Corners [NC1 and NC2], two Near
217 Edges [NE], two Far Corners [FC] and one Far Edge [FE]). This allows us to create two types of
218 fragment scenes. The first type pairs one of the fragments in the surround with the CE ('with-CE'). Two
219 additional scenes contain respectively only the CE, or all the fragments. The second type are identical
220 scenes but without the CE ('without-CE'). For 'with-CE' scenes we defined the surround influence of a
221 fragment as the unit's response to the combination of that fragment and the CE, subtracted by the
222 response to the CE-only scene (see Methods). For 'without-CE' scenes, the surround influence of a
223 fragment was defined as the response to a scene with that fragment, subtracted by the response to a full-
224 gray scene. Figure 2E displays the data for one example BOS unit in PredNet. First, we noticed that the
225 absolute value of surround influence in 'with-CE' scenes is larger than in 'without-CE' scenes. This was
226 the case for each PredNet module with at least 10 BOS units (Fig. 2F for E_2 and other modules in SI Fig.
227 5). This is similar to BOS neurons in the visual cortex (Figs. 2G,H)³¹. Second, we compared the
228 modulation effect between fragments on the preferred side and the non-preferred side. The preferred and
229 non-preferred sides were determined solely from the responses to standard square scenes (Fig. 1A).
230 Despite this, we found that for all modules with more than 10 BOS units, the surround influence for
231 most fragments is significantly more negative when they are presented on the non-preferred side
232 compared to the preferred side (Fig. 2F; SI Fig. 5). This is similar to BOS neurons in the visual cortex
233 (Figs. 2G, H). These data indicate that BOS tuning in PredNet does not result from a single hotspot in
234 the surround, but that multiple fragments collectively contribute, as is the case for BOS neurons in the
235 brain³¹.

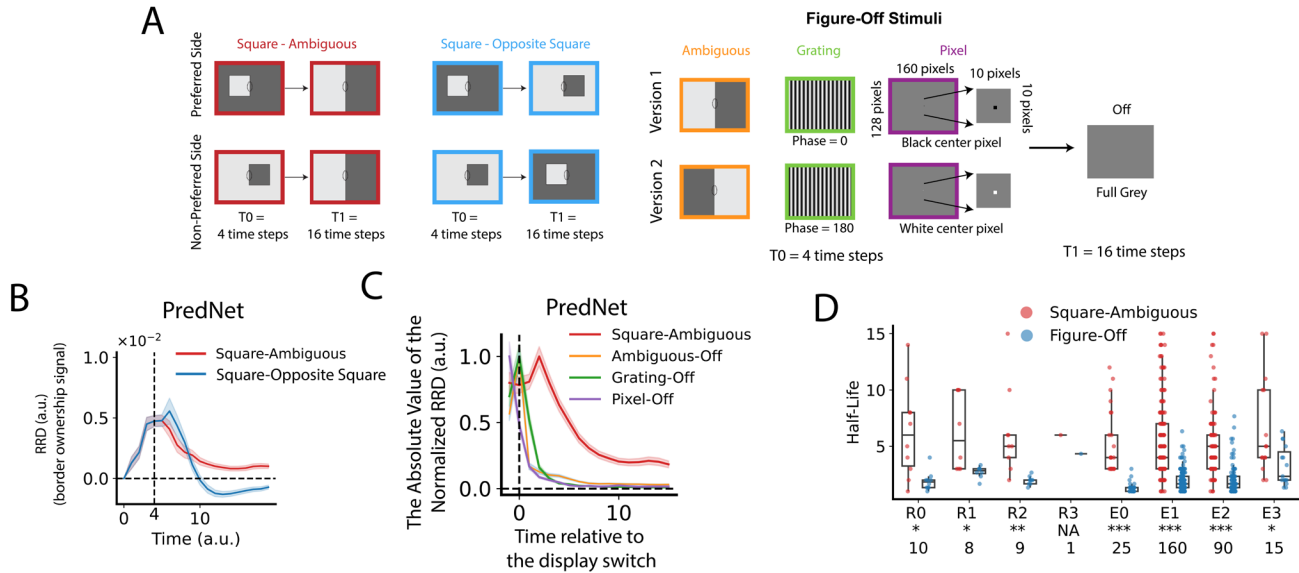
236 **PredNet's BOS units exhibit hysteresis, similar to BOS neurons in the brain**

237 A remarkable characteristic of BOS neurons in the brain is that the BOS signal persists for hundreds of
238 milliseconds, even when the contextual information that defines the side of ownership disappears^{24,25}.
239 We tested if BOS units in PredNet also exhibit this phenomenon. We used a Square-Ambiguous
240 sequence similar to what was used in physiology experiments²⁴. The sequence consists of a full square
241 scene (Figure 1A) in the first four time steps, which transitions into a scene with a border that is
242 ambiguous for border ownership (Figure 3A, left). We presented these sequences to PredNet and
243 computed the time course of the relative response difference (RRD), defined as the difference in
244 response between the preferred and non-preferred square sides, normalized by the average response (see
245 Methods). The RRD for BOS units remains positive for multiple time steps (Fig. 3B left, red function;
246 SI Figure 6). The units thus respond differently to the ambiguous scene (which is identical in the two
247 sequences), depending on stimulus history, a phenomenon called hysteresis. BOS neurons in macaque
248 visual cortex show a similar hysteresis (compare with Fig. 2A in O'Herron and von der Heydt, 2009²⁴).

249 To determine whether this persistent BOS signal is longer than the typical signal decay, we analyzed the
250 response for two control sequences. The first is Square-Opposite Square, which starts with a full square
251 and then switches to another full square image with opposite border ownership (and opposite luminance,
252 so that the contrast polarity of the central border remains the same; Fig. 3A, middle). The RRD for this
253 sequence decays much faster and stabilizes at a negative value, reflecting the switch in BOS (Fig. 3B,
254 left, blue function; SI Figure 6). Again the same pattern occurs in BOS neurons in the brain (compare
255 with Fig. 2A in O'Herron and von der Heydt, 2009²⁴). The second control sequence is Figure-Off, in
256 which a simple scene (three subtypes: ambiguous, grating or pixel) is followed by a full gray scene (Fig.
257 3A, right). Again the RRD of PredNet's BOS units decays faster to these sequences than to the Square-

258 Ambiguous sequence (Figs. 3C,D, SI Figure 6B). Across all modules with at least 10 BOS units the
 259 RRD half-life is significantly longer for Square-Ambiguous sequences than for Figure-Off sequences
 260 (Fig. 3D). Together, we find that BOS signals in PredNet have similar dynamic characteristics as BOS
 261 neurons in the brain: the BOS signal persists when contextual information disappears such that the side
 262 of BOS becomes ambiguous, but quickly updates when the context indicates a switch in BOS.

263



264

265 **Figure 3: BOS signals in PredNet exhibit hysteresis, similar to BOS neurons in the brain.** (A) Three
 266 sequences with scene changes were used: squares transitioning to ambiguous borders (Square-Ambiguous),
 267 squares transitioning to squares with opposite border ownership (Square-Opposite Square), simple scenes
 268 (Ambiguous, Grating, or Pixel) transitioning to a full gray scene (Figure-Off). Stimuli were presented at the
 269 orientation for which |BOI| was maximal. (B) The relative response difference (RRD, see Methods) represents the
 270 difference in response between scene sequences that start with a square on the preferred and the non-preferred
 271 side (for Square-Ambiguous or Square-Opposite Square), or between version 1 and version 2 (for Figure-Off).
 272 Panel shows RRD of BOS units from PredNet (E_2 module, $n = 132$ units). Functions plotted in the same format as
 273 Fig. 2A in O’Herron and von der Heydt, 2009²⁴. Line and error bands represent the mean and SEM. (C) Mean and
 274 SEM of the absolute value of the normalized RRD (normalized to maximal value) across BOS units in the E_2
 275 module for different sequences. (D) Half-life is defined as the number of time steps after which RRD is reduced
 276 to half of its maximum. Each dot corresponds to one BOS unit. Figure-Off data shows the average across the three
 277 subtypes shown in A. Only units for which half-life was defined for all conditions were included in this panel (see
 278 Methods). Numbers at the bottom indicate the number of included units per module. Asterisks indicate the
 279 statistical significance of the difference in half-life between Square-Ambiguous and Figure-Off: NA: not
 280 applicable; *p<0.05, **p<0.01, ***p<0.001 (Wilcoxon signed-rank test).

281

282 BOS units contribute more to prediction than non-BOS units for videos with moving objects

283 Our data presented thus far demonstrate that units with brain-like tuning for BOS exist in PredNet, a
 284 network trained to predict future visual input in video sequences. This suggests that BOS units
 285 specifically aid in predicting future video frames. To test that, we conducted ablation experiments in
 286 PredNet. We presented Translating-Square videos (40 unique videos in which a square moves at a
 287 constant velocity, SI Fig. 7, top) to PredNet. We measured the prediction performance of PredNet to

288 these videos, both before and after ablating either BOS units or non-BOS units (i.e. candidate units that
289 did not pass the criterion for BOS-selectivity, see Methods).

290 The impact of unit ablation on video prediction is shown in SI Figure 9 (top row). Here, we introduce
291 the metric “relative prediction mean squared error (RPE),” defined as the normalized difference (post-
292 vs. pre-ablation) of the mean squared prediction error (see Methods). A positive RPE represents an
293 increase in prediction error after ablation. To quantify the overall effect of ablation in each module, we
294 measured the slope of the relation between RPE and number of ablated units using linear regression, and
295 a bootstrapping test to assess the statistical significance of this slope between ablating BOS units or non-
296 BOS units (indicated with red symbols in SI Fig. 9, top row). We find that the RPE is significantly
297 higher when BOS units were ablated than when non-BOS units were ablated for most modules. We
298 wondered if this could be explained by a difference in responsiveness: BOS units may respond more to
299 these video frames than non-BOS units. To explore that possibility, we subsampled the populations to
300 ensure there were no statistically significant differences in response magnitude to the videos (Wilcoxon
301 rank-sum test, $p > 0.5$, see SI Figure 8 and Methods). The ablation experiment in these subsampled
302 populations shows the same pattern, ruling out that the RPE difference is due to a difference in average
303 response (Fig. 4, top row). The data thus indicate that BOS units contribute more than non-BOS units in
304 predicting future frames for these videos.

305 We wondered if BOS units also contribute to prediction of videos with multiple objects. We generated
306 videos with several squares that were randomly positioned, and moved in random directions (SI Fig. 7,
307 middle). When we performed the same ablation experiment for these videos, we find the same pattern:
308 BOS units typically contribute more to prediction than non-BOS units, even though, again, PredNet was
309 not exposed to such videos during training (Fig. 4 bottom, SI Fig. 9 middle).

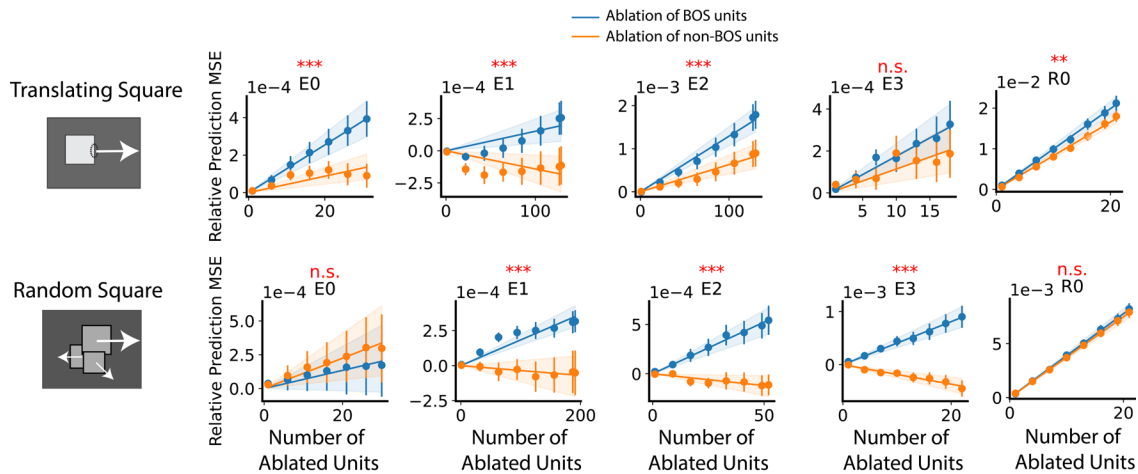
310 Finally, we wondered if BOS units aid in prediction with any video. We performed the same experiment
311 in a set of 41 natural videos from the KITTI database²⁸ (SI Fig. 7, bottom). This is the same database
312 that was used to train PredNet, but we only included videos that were not used during training. BOS
313 units in the E_2 contribute more to prediction than non-BOS units (SI Fig. 9). Note that these videos are
314 much higher-dimensional than the translating square and random square videos, and there is a high
315 degree of heterogeneity within the small set of 41 videos. This results in smaller overall RPEs when
316 averaged across videos than for the square videos, and not enough statistical power to precisely estimate
317 RPE in the subsets of units with similar responsiveness (SI Fig. 10).

318 Together these experiments suggest that BOS units emerge in PredNet because they contribute more to
319 prediction than non-BOS units for videos with moving objects.

320

321

322



323
324

325 **Figure 4. Ablating BOS units in PredNet increases prediction error more than ablating non-BOS units for**
 326 **videos with moving objects.** The left shows an example frame of each video type (arrows indicate motion and
 327 are not part of the frame). Translating Square videos show a square moving at constant speed; Random Square
 328 videos show a random number of squares of different sizes, initialized at random positions and moving at random,
 329 constant velocities (see also SI Figure 7). Right panel shows the relative prediction mean squared error (RPE) for
 330 different numbers of ablated units. RPE measures the relative change of prediction error due to ablation. Non-
 331 BOS units are candidate units that do not pass the criterion of BOS selectivity. Dots and error bars denote
 332 respectively the mean and SEM of the RPE across 10 randomly chosen video samples. The RPE of one video
 333 sample is the average RPE of 10 samples of unit ablation (Methods). The solid line indicates the best linear fit,
 334 with bands indicating the 95% confidence interval. The red text above the panels indicates whether the slopes of
 335 the lines differed significantly between BOS- and non-BOS-unit ablation. n.s.: not significant; **: $p < 0.01$; ***: p
 336 < 0.001 (bootstrapping test). Modules R_1 , R_2 and R_3 contain a small number of candidate units and are therefore
 337 not included in this analysis.

338
339
340
341

Discussion

The assignment of borders to foreground surfaces is thought to be a key step in visual scene segmentation^{13,32}, and a substantial fraction of neurons in visual areas V2 and V4 of the primate brain signal this ownership of local borders¹⁻³. It is poorly understood why the brain resorts to this particular representation. Here we discovered that units selective for BOS also emerge in an artificial neural network, PredNet²⁶, that was trained to predict future input in natural videos. Importantly, the network was not explicitly trained to distinguish foreground from background or to identify objects in visual scenes. Interestingly, BOS units in PredNet share several properties with BOS neurons in the brain (robustness for different positions, orientations, and sizes^{1,3}; asymmetric functional effects of object fragments on opposite sides of the border³¹; BOS hysteresis²⁴), suggesting that these signals are functionally similar to those in the brain. Finally, we found that ablation of BOS units affects prediction accuracy more than ablation of non-BOS units. Overall, our results suggest that BOS units might emerge in neural networks trained on natural, complex dynamic input primarily because they are particularly helpful to efficiently process such input, even if segmentation is not required.

PredNet's architecture was inspired by the predictive coding framework. This theory proposed that a major function of the sensory cortex is to predict incoming sensory stimuli³³⁻³⁷. The hierarchical organization of visual cortical areas is proposed to compute an internal model of the external world, and feedback from areas higher in the hierarchy (e.g. V4, IT) is thought to reflect predictions from this internal model, which is then compared with incoming sensory stimuli in lower areas (e.g. V1, V2)³³. There are hints suggesting that how brain circuits compute border ownership may be understood in this framework^{13,38-40}. For example, area V4 has been proposed to contain grouping cells which compute proto-object representations with short latency, i.e. an early prediction of the shape and location of objects in the scene. Feedback from such cells could explain border ownership signals in lower areas^{6,13,31,41}, and a recent study indeed found evidence that supports the existence of grouping cells in V4²³. Response dynamics and laminar organization of BOS neurons align better with feedback models than with alternatives that solely rely on intra-areal horizontal connections or feedforward connections^{3,4,31,42-46}. Moreover, the phenomenon of BOS hysteresis indicates that BOS neurons persistently signal the most likely scene organization even if contextual information disappears, but quickly update when sensory information inconsistent with the current internal model appears^{24,47}. The present data provide complementary evidence that there may indeed be a tight link between predictive coding and how neural networks compute BOS. We showed that a predictive coding inspired architecture can lead to BOS signals, with properties very similar to those in the brain, even without explicitly training the network to localize or identify objects in visual scenes.

A prior study showed that PredNet units signal illusory contours and end-stopping¹⁷. The emergence of BOS signals as well as these other extra-cRF phenomena under the predictive coding framework raises a question: do these phenomena result from a single hierarchical neural computation? Several lines of prior research are consistent with that possibility^{17,30,33,34,36,48}. A complete answer to this question is hard to obtain by solely doing physiology experiments: detailed maps of neural connections are often unavailable, and it is challenging to precisely manipulate these connections. ANNs have the unique advantage of possessing complete connection profiles^{22,49,50}, and allow one to perform ablation studies. Our work thus establishes PredNet as a useful complementary tool towards achieving an understanding of these computations.

385 PredNet's E modules have been interpreted as being akin to superficial layers (L1/2/3), and the R
386 modules as akin to deeper layers (L5/6) of the visual cortex, following the proposed functional
387 specialization of cortical layers in predictive coding^{17,27,33}. The presence of BOS signals in both E and R
388 modules aligns with physiology, where BOS neurons exist in both superficial and deeper layers³.
389 However, one should be cautious to equate E and R modules to different cortical layers. For example,
390 the R modules have lateral connections, which the E modules lack these, unlike in physiology where
391 lateral connections exist in both superficial and deep compartments^{36,51,52}. Further studies are needed to
392 understand the functional role of different areas and layers in this hierarchical computation and the
393 communication between them. Because of the flexibility to manipulate network architecture and
394 connections, ANNs are a useful complementary tool in such studies^{17,27}.

395 Our discovery of BOS units in PredNet and the ablation experiments indicate that BOS neurons may be
396 useful for video prediction. To predict future visual input, it is useful to predict object motion⁵³. Objects
397 typically move as a whole, i.e. pixels within object boundaries most likely move together⁵⁴. Because
398 BOS units indicate which pixels belong to an object surface, they may help to predict by allowing the
399 system to easily apply a uniform optical flow to objects. Indeed, in computer science, incorporating
400 optical flow^{55,56}, disentangling object motion from content⁵⁷⁻⁶⁰, and separating foreground objects from
401 background^{54,61} have been shown to improve video prediction performance. Beyond video prediction, in
402 object recognition, deep neural networks have been criticized for relying mostly on textural information
403 to recognize object categories rather than on object shapes⁶²⁻⁶⁴, in contrast to human visual
404 perception^{65,66} (but perhaps more akin to mouse visual perception⁶⁷). Explicitly embedding a BOS unit
405 module may guide neural networks to rely more on shapes, and potentially achieve more robust
406 recognition as well as prediction.

407 Overall, our work demonstrates that brain-like BOS signals emerge in a self-supervised network trained
408 to predict future input. This implies a shift from the traditional view of BOS as a static 'what stream'
409 operation towards a computation that is highly beneficial to predict future input in natural dynamic
410 environments.

411

412

413

414

Methods

PredNet architecture

415
416 In this study, we utilized the artificial neural network PredNet, which was developed and trained by
417 Lotter et al. (2017)²⁶ (code is available at: <https://github.com/coxlab/prednet>). Here we briefly
418 summarize PredNet's architecture and how it was trained. PredNet is an artificial neural network (ANN)
419 that has four layers (labeled as 'l'). Each layer consists of four types of modules: the Representation
420 module (R_l), the Prediction module (\hat{A}_l), the Prediction target module (A_l), and the Prediction Error
421 module (E_l). Updating unit activities in PredNet involves two main stages at each time step:

422 Top-to-Bottom Update: The network updates the R modules from top to bottom at each time step. Each
423 R_l module gets inputs from the R_{l+1} module and the E_l module. This updating process goes from the R_3
424 module to the R_0 module in sequence. The R_0 module then generates a predicted current video frame
425 (\hat{A}_0).

426 Bottom-to-Top Update and Error Calculation: The update process then reverses, proceeding from
 427 bottom to top. The network calculates the prediction error by comparing \hat{A}_0 with the actual next video
 428 frame, A_0 . This error is bifurcated into positive and negative parts (akin to biological ON-center and
 429 OFF-center neurons). Positive and negative errors are grouped in the E_0 module. E_0 then outputs a target
 430 prediction A_1 , which gets compared with \hat{A}_1 produced from R_1 . The error from this comparison is the E_1
 431 module. The network continues this process up to the final layer (layer 3).

432 Mathematically, the PredNet dynamics are defined by

$$\begin{aligned}
 A_l^t &= \begin{cases} x_t, & \text{if } l = 0 \\ \text{MAXPOOL}(\text{RELU}(\text{CONV}(E_{l-1}^t))), & l > 0 \end{cases} \\
 \hat{A}_l^t &= \text{RELU}(\text{CONV}(R_l^t)) \\
 E_l^t &= [\text{RELU}(A_l^t - \hat{A}_l^t); \text{RELU}(\hat{A}_l^t - A_l^t)] \\
 R_l^t &= \text{CONVLSTM}(E_l^{t-1}, R_l^{t-1}, \text{UPSAMPLE}(R_{l+1}^t))
 \end{aligned} \tag{1}$$

433 where t is the time step, x_t is the actual video frame. ConvLSTM uses a tanh activation function, which
 434 means that the R module activation can be negative (possible values range from -1 to 1). Because
 435 biological neurons do not have negative spike rate, PredNet unit's response was defined in this study as
 436 the unit activation plus one, i.e. the response baseline was shifted by +1 in all modules (after PredNet's
 437 computation was completed, thus this did not affect PredNet's algorithm). The PredNet architecture
 438 contains 3, 48, 96, and 192 convolution channels in layers 0 to 3, respectively. The input image size is
 439 128 by 160 pixels. The number of units in R modules are respectively 61,440 in R_0 , 245,760 in R_1 ,
 440 122,880 in R_2 , and 61,440 in R_3 . The A_l and \hat{A}_l modules have the same number of units as the R_l
 441 module. Due to the bifurcation of positive and negative error, E modules have twice the number of units
 442 compared to the R modules.

443 The training loss function is applied on the prediction error

$$L_{train} = \frac{1}{N_t N_0} \sum_t \sum_{n_0} E_0^t \tag{2}$$

444 where N_t is the number of time steps used in training, N_0 is the number of E_0 units. The training utilized
 445 the KITTI dataset, which contains videos recorded from car-mounted cameras in Germany. Videos were
 446 segmented into sequences of 10 continuous frames. These frames were then center-cropped and
 447 downsampled to a resolution of 128 by 160 pixels. The parameters of PredNet were optimized using
 448 backpropagation with the Adam optimizer.

449 Classical Receptive Field of Units

450 We measured the classical receptive field (cRF) of units in PredNet using sparse noise stimuli (SI Figure
 451 1), similar to the approach used in physiology. We created an image (128 x 160 pixels) with one pixel
 452 set to either white or black, while all others were set to gray (gray level = 0.5, scales from 0 to 1). 40,960
 453 (128 × 160 × 2 where the factor two is for black and white pixel) unique images (128 x 160) were
 454 generated, each featuring a distinct single pixel, in either white or black. These images were repeated for
 455 four time steps, yielding a total of 40,960 sequences. For each unit, recorded activity to these sequences
 456 was summarized into two heatmaps (each size 128 × 160), each representing responses to respectively
 457 white-pixel and black-pixel scenes. For example, the white heatmap's i, j entry is the single unit's time-
 458 averaged response to a scene with a white pixel located at i, j (gray otherwise).

459 The two heatmaps (for one unit) were then z-scored and converted to absolute values. These heatmaps
460 were merged into one heatmap by taking the maximum absolute values for each entry. This merged
461 heatmap summarizes the unit's maximum response to a pixel at each location irrespective of its color
462 (white or black). The cRF for each unit was defined as the union of the pixel positions for which the
463 absolute values of the maximal z-scores across both heatmaps exceed 1.

464 **Standard Square Stimuli**

465 Scenes with square objects are commonly used in neurophysiological studies to assess whether a unit is
466 selective for BOS^{1-3,23} and this selectivity is known to extend to natural images⁵. We used similar
467 scenes, consisting of a square with a size (width) of 50 pixels, positioned with one border centered at the
468 center of the scene (central border). The color of the square and the background can be either light (gray
469 level = 0.33 on a scale from 0 to 1) or dark gray (gray level = 0.66), but they are always different from
470 each other in a given scene. These square scenes can be defined mathematically by three parameters.
471 The first parameter, α , denotes the square's orientation, with a range from 0 to 180 degrees. The second
472 parameter, β , is a binary variable indicating which side the square is given a fixed orientation (i.e. side
473 of ownership). The final parameter, γ , is a binary variable that indicates the contrast polarity across the
474 central border. For each square orientation defined by α , there are four possible square scenes,
475 determined by different combinations of β and γ . Each of these scenes is repeated over 20 time steps.
476 We used 10 different orientations (equally spaced by 18°).

477 **Candidate Unit Selection**

478 To define selectivity for border ownership, it is important to verify that the units under examination
479 respond to changes in border ownership rather than to low-level stimulus changes within the cRF.
480 Therefore, similar as in neurophysiology studies, we restricted our analysis to units that passed the
481 following two criteria (termed 'candidate units'). First, the unit's cRF must include the center of the
482 scene. Because the central border of the square scenes was placed exactly in the scene center, this
483 ensured that the unit's cRF includes the center of this border. Second, the cRF must fit within a circle
484 centered at the center of the scene and with a radius of 20 pixels. Because the square size (width) is 50
485 pixels, this makes sure that the cRF does not overlap with any other border of the square besides the
486 central border.

487 **Averaged Border Ownership Index across orientations (B_{av})**

488 Similar to neurophysiology studies¹⁻³ we quantified tuning for border ownership using the Border
489 Ownership Index (BOI). This is computed from the response of PredNet units to standard square scenes.
490 The BOI is defined as

$$BOI(\alpha) = 2 \times \frac{Res(\alpha, 1,0) - Res(\alpha, 0,0) + Res(\alpha, 1,1) - Res(\alpha, 0,1)}{Res(\alpha, 1,0) + Res(\alpha, 0,0) + Res(\alpha, 1,1) + Res(\alpha, 0,1)} \quad (3)$$

491 where $Res(\alpha, \beta, \gamma)$ is the unit's time-averaged (between 0 and 19 time steps) responses to a square scene
492 specified by orientation α , side-of-ownership β and contrast polarity γ . The sign of the BOI thus
493 indicates which side (β) of BOS (for a given orientation) the unit prefers, and the magnitude indicates
494 the strength of the BOS tuning.

495 To evaluate the overall BOS selectivity across orientations, we defined B_{av} as the circular average of the
496 BOI across α . Similar to BOI, the magnitude of B_{av} is a measure of the strength of BOS tuning, and its
497 angle indicates the unit's preferred side of BO.

498 We evaluated the statistical significance of B_{av} using a permutation test. In this test, we shuffled the
499 labels that signified the side of BOS (β) for each orientation α . These data were then used to compute a
500 shuffled BOI(α) and B_{av} . This procedure was repeated 5,000 times to generate a set of 5,000 B_{av} values
501 after shuffling, for each unit. Denoting the quantile of the unshuffled B_{av} among the shuffled B_{av} as Q ,
502 the p-value (two-tailed) was estimated as $2 \times \min\{Q, 1 - Q\}$. Units with a p-value less than 0.05 were
503 defined as BOS units. 95% confidence intervals on proportions of units for which B_{av} was significant
504 were computed using Wilson score⁶⁸.

505 Note that the values of B_{av} and BOI reported here cannot easily be compared with similar indices in
506 neurophysiology, because these values change when the DC level of unit activity is changed. As
507 mentioned above, to avoid negative values for unit activity in PredNet, we arbitrarily increased activity
508 levels by +1. Furthermore, the average BOI across time depends on when the response starts relative to
509 the duration of the analysis window. This is at ~50% of the window duration for the unit shown in Fig.
510 1D (top), whereas in physiology studies this is typically closer to ~10%. For example, the activity
511 functions shown in Figure 1D (top panel) show a BOI of 0.0149 at time step 10, but computing this
512 without adjusting the unit activation (i.e. without +1) leads to BOI = 0.68. Zhou et al. use 'response
513 ratio' to quantify the magnitude of BOS tuning, defined as the ratio of the mean response to non-
514 preferred BOS over the mean response to preferred BOS. For the activity functions shown in Fig. 1D
515 (top panel) this value is 0.561 (averaged across analysis window), well within the range of values found
516 for neurons in the macaque visual cortex¹.

517

518 Analysis of BOS Unit Responses to Different Square Orientations, Positions, and Sizes

519 In these experiments, varied parameters were square orientation (α), side-of-ownership (β), contrast
520 polarity (γ), position along the orientation (d), and size (s). We first measured the response to a set of
521 four standard square scenes (Figure 1A). For each unit, the orientation α is fixed at the orientation with
522 the maximum absolute BOI. The position is zero, indicating that the square border intersects exactly
523 with the scene center, and the square size (width) is 50 pixels. BOS units' responses were averaged over
524 time and contrast polarity. The β value with the larger averaged unit response was defined as the
525 preferred side (β_p), whereas the opposite was defined as the non-preferred side (β_{np}). These preferences
526 were solely determined by the standard square scenes.

527 We then examined the effect of changing square size. All other parameters remained the same as in the
528 standard square scenes stated above, except for square size. Eight square sizes were used, ranging from
529 10 to 80 pixels. For each unit i and each square size s_j , we computed the responses averaged across time
530 and contrast polarity, yielding $\bar{r}_{i,j}(\beta_p)$, $\bar{r}_{i,j}(\beta_{np})$. We then normalized two response arrays of each unit
531 i : $\tilde{r}_{i,j}(\beta) = \bar{r}_{i,j}(\beta) / \sum_j \bar{r}_{i,j}(\beta)$, where β can be β_p or β_{np} . Figure 2D (left panel) displays the time
532 course of $\tilde{r}_{i,j}(\beta)$ across units i . For each unit i and square size s_j , we computed a BOI as the difference
533 in response between the β_p and β_{np} , i.e. $BOI_{i,j} = 2 \times (\tilde{r}_{i,j}(\beta_p) - \tilde{r}_{i,j}(\beta_{np})) / (\tilde{r}_{i,j}(\beta_p) + \tilde{r}_{i,j}(\beta_{np}))$. We
534 performed a bootstrapping test to assess statistical significance of this metric. A BOI dataset consisted of
535 $D^{BOI} = \{BOI_{i,j}\}$ for all units i and square sizes s_j . We obtained 10,000 bootstrap samples D^S from this

536 dataset. For each D^S , we computed an averaged BOI, denoted as BOI^S . The p-value was estimated as
537 $p = 1 - Q$, where Q is defined as the quantile of 0 among all BOI^S . If p-value was smaller than 0.05,
538 we concluded that the BOI averaged across size was statistically significantly positive in the population
539 of BOS units (SI Figure 4).

540 The same procedures apply to varying square position, simply replacing square size with square position
541 (SI Figure 4). Fifteen square positions were used, ranging from -30 to 26 pixels.

542 When examining the unit's response to different orientations, we created square scenes with 10 possible
543 orientations (equally spaced between 0 and 180 degrees), keeping the position at 0 and size at 50 pixels.
544 Units' responses were collected to compute the BOI for each orientation using equation (3). Data from
545 one example unit is shown in Figure 2A (left panel). To evaluate the statistical significance of BOI for a
546 given orientation, we compared the unshuffled BOI to that in a null distribution. Unlike biological
547 neurons, which differ in response from trial to trial, PredNet does not have noise. In order to obtain
548 sufficient data to generate a shuffled distribution, for each orientation, we varied square size (10
549 different sizes were considered, conceptually mimicking 10 “repeated trials”). The unshuffled BOI for a
550 given orientation was computed for this orientation across square size. The null distribution for BOI
551 distribution was obtained by shuffling the labels indicating the side-of-ownership β (i.e., border
552 ownership), separately for each square size and contrast polarity (5,000 shuffles). The quantile (Q) of the
553 unshuffled BOI within shuffled BOI set was computed. The p-value (two-tailed) was estimated as
554 $2 \times \min\{Q, 1 - Q\}$. If the p-value is less than 0.05, BOI along an orientation was said to be statistically
555 significant (indicated in Fig. 2A as filled circles). The above procedure resulted a subset of orientations
556 with statistically significant BOIs. The span of each BOS unit was computed as the difference between
557 the two most distant preferred object locations (circular distance between the two angles corresponding
558 to those locations). 95% confidence intervals on proportions of units for the span was smaller than a
559 certain value were computed using Wilson score⁶⁸.

560 **Square Fragment Stimuli**

561 The squares in the square scenes can be divided into eight fragments³¹: the Central Edge (CE), which is
562 the one in the middle of the scene; there are two Near Corners (NC), two Near Edges (NE), two Far
563 Corners (FC), and one Far Edge (FE). To examine how these fragments modulate the activity of BOS
564 units, the four standard square scenes (Figure 1A, orientation aligns with the preferred BOS orientation
565 for each unit) were converted into fragmented square scenes, as described below.

566 To isolate one fragment, a 2D Gaussian filter ($\sigma = 5$ pixels) was applied at the fragment's center. This
567 kept the fragment's central region largely unaltered, while the parts of the scene further away gradually
568 fade to a uniform gray (gray level = 0.5 on a scale from 0 to 1). For scenes with multiple fragments (e.g.
569 ‘All’), a Gaussian filter was applied to each fragment. Note that the smallest distance between two
570 fragment centers is 25 pixels, thus much larger than σ , resulting in negligible interference between
571 filtered fragments at different locations.

572 Using this Gaussian filter method, we created 9 scenes with a Central Edge (‘with-CE’) and 9 scenes
573 without a Central Edge (‘without-CE’) for each of the four standard square stimuli (Figure 1A). Among
574 the with-CE scenes, one scene only has the CE fragment, seven scenes have the CE and one additional
575 fragment, and one scene has all fragments. The without-CE scenes are similar to the with-CE scenes,
576 except that they do not contain the CE fragment. Thus the “all fragments” without-CE scene contains 7
577 fragments. Each scene is presented during 20 time steps.

578 **Processing of Units' Responses to the Square Fragment Stimuli**

579 The NC fragment could potentially partially intersect with the cRF. To prevent this, we limited this
580 analysis to the subset of BOS units whose cRFs fitted within a circle of 30 pixels diameter centered at
581 the center of the scene. This more conservative selection yielded 30, 145, 71, 5 units from respectively
582 E0 to E3, and 1, 3, 2, 0 units from respectively R0 to R3. For with-CE scenes, the surround influence of
583 square fragment X is defined as the unit's response to the X + CE scene subtracted by the response to
584 the CE scene. Similarly, for without-CE scenes, the surround influence of square fragment X is defined
585 as response to the X scene subtracted by a full gray scene. If X is FE, the surround influence of X is
586 computed as above. Otherwise (X = FC, NE, NC), the surround influence of X is the average of the
587 surround influences of two conjugate edges (e.g., CE1 and CE2).

588 The surround influences of X for all BOS units were computed, resulting in a list where the length
589 equals the number of BOS units. To avoid bias in mean estimation due to outliers, outliers (1.5
590 interquartile range below the first quantile or above third quantile) were removed before computing the
591 sample mean and SEM (Fig. 2F and SI Fig. 5). However, all units were included when performing
592 statistical tests (indicated by figure caption).

593 **Square-Ambiguous, Square-Opposite Square, and Figure-Off Sequences**

594 Each trial in the Square-Ambiguous sequences consisted of 20 time steps, broken down into two phases.
595 Initially, Scene 0, one of the four standard square scenes (Figure 1A), was displayed during four time
596 steps ($T_0 = 4$). Subsequently, Scene 1 was shown during 16 time steps ($T_1 = 16$). Scene 1 only
597 contained a central border that divides the whole image into a left and a right half; hence the side of
598 ownership of this border was ambiguous. The contrast polarity and orientation of Scene 1 were
599 consistent with Scene 0 (i.e. the information in the cRF was the same).

600 Similarly, the Square-Opposite Square sequences started with one of the four standard square scenes as
601 Scene 0. Scene 1 was a version of the square scene with reversed BOS, but maintaining contrast polarity
602 for the central border. For example, if Scene 0 was panel 1 in Figure 1A, then Scene 1 was panel 2 in
603 Figure 1A.

604 For Figure-Off Sequences, Scene 1 was always a full gray. Scene 0 depended on the subtypes:
605 Ambiguous-Off, Grating-Off, and Pixel-Off sequences. For Ambiguous-off, Scene 0 was an ambiguous
606 border. It had two versions that vary in contrast polarity. In Grating-Off sequences, Scene 0 was a
607 grating with a 10-pixel spatial period, and it had two versions with grating phases of either 0 or 180
608 degrees. For Pixel-Off sequences, Scene 0 was gray except for a single pixel at the center, which was
609 either white or black corresponding to two versions.

610 All scenes were generated such that the orientation corresponds to that for which each unit's |BOI| was
611 maximal.

612 **Relative Response Difference**

613 The Relative Response Difference (RRD, used in the result section "PredNet's BOS units exhibit
614 hysteresis, similar to BOS neurons in the brain" and Figure 3) is $(a - b)/(a + b)$, where a indicates the
615 time-averaged response to preferred stimuli, and b indicates the time-averaged response to non-preferred
616 stimuli. Which stimulus was preferred only depended on the averaged response to Scene 0.

617 RRD half-life was defined as the earliest time after the scene switch where the absolute value of RRD
618 was less than half of its maximum. The half-life across the three types of Figure-Off sequences were
619 averaged in Figure 3D. For this analysis, we only included units for which the half-life of all three types
620 of Figure-Off sequences could be measured (exclude RRD that never dropped to half of its maximum
621 within the analysis window). This yielded 10 out of 22, 8 out of 25, 9 out of 12 and 1 out of 2 BOS units
622 in respectively R0, R1, R2 and R3; and 25 out of 32, 160 out of 199, 90 out of 131, and 15 out of 22
623 BOS units in respectively E0, E1, E2 and E3. The Wilcoxon signed-rank test was used to compare half-
624 life between Square-Ambiguous sequences and Figure-Off sequences.

625 **Three video types for Ablation Experiment**

626 We generated three types of videos to evaluate PredNet's prediction performance (examples shown in SI
627 Figure 7). (1) Translating Square videos include a square that moves at a constant speed and direction.
628 Square size is 50 pixels and oriented such that the central border had a vertical orientation (square gray
629 level = 0.33 and background gray level = 0.66 on a scale from 0 to 1). The initial position and velocity
630 of the square were chosen such that the square was always in the scene center in the 10th frame. Forty
631 translating square videos were created, corresponding to 40 evenly spaced moving directions (equally
632 spaced between 0 and 360 degrees). (2) Random Square videos: each of these videos featured a random
633 number of squares (between 1 and 5). At the beginning of each video, each square's central position was
634 randomly set in the scene. The size of each square was also randomly chosen (between 10 and 50
635 pixels), and the x and y components of each square's velocity were randomly set at a value between -2
636 and 2 pixels/frame. Forty random videos were generated. (3) KITTI testing videos: 41 videos from car-
637 mounted cameras were used, which were not used during PredNet's training²⁶. For all video types, each
638 video consisted of 20 frames.

639 **Subsampling BOS and Non-BOS Units to Reduce Their Response Differences**

640 Unit activity in response to the videos were squared and averaged across all videos and time steps for
641 each video type, resulting in Mean Squared Response (MSR). For each module and video type, we have
642 two sets of MSR, one for the BOS units and another for the non-BOS units, denoted as $D^{bos} =$
643 $\{r_0^{bos}, r_1^{bos}, \dots, r_n^{bos}\}$ and $D^{non-bos} = \{r_0^{non-bos}, r_1^{non-bos}, \dots, r_m^{non-bos}\}$, respectively, where n and m
644 representing the number of BOS and non-BOS units in one module.

645 For each of the D^{bos} and $D^{non-bos}$, we subsampled $k = \min\{n, m\}$ units (1,000 samples). This resulted
646 in 1,000 pairs of sampled datasets, denoted as D_s^{bos} and $D_s^{non-bos}$, with s ranging from 1 to 1,000. For
647 each pair, we computed a score to measure the similarity between datasets in a pair

$$\phi_s = [\text{mean}(D_s^{bos}) - \text{mean}(D_s^{non-bos})]^2 + [\text{median}(D_s^{bos}) - \text{median}(D_s^{non-bos})]^2 \quad (4)$$

648 where the $\text{mean}(\cdot)$ and $\text{median}(\cdot)$ represent those quantities of the dataset. The dataset pair with
649 smallest score ϕ_s was subjected for further statistical analysis, using the Wilcoxon rank-sum test and the
650 t-test. If both p values were larger than 0.5, we considered the dataset pair as our final subsampled
651 datasets. If not, we reduced k by 1 and repeated the procedure above. This whole procedure makes sure
652 that both BOS and non-BOS populations have the same number of units (equal to k), and their MSRs do
653 not show significant difference. SI Figure 8 displays the MSR of the obtained subsampled unit
654 populations.

655 Compute the Prediction Error of the Ablation Experiment

656 For each video type, we created N_α bootstrapped samples, each containing N_α videos. We denoted v_a^α as
 657 the a^{th} video in the α^{th} bootstrapped sample, with α ranging from 0 to $N_\alpha - 1$, and a from 0 to $N_\alpha - 1$.
 658 In this study, $N_\alpha = N_a = 10$.

659 For each video v_a^α , we performed the ablation experiment several times, for different samples of ablated
 660 units, in each module separately. We varied the number of ablated units n (ranges from 1 to
 661 $\min\{N_{bos}, N_{non-bos}\}$, where N_{bos} and $N_{non-bos}$ indicate respectively the number of BOS and non-BOS
 662 units available in the module). For each n , we generated $N_u = 10$ bootstrapped unit samples from the
 663 unit pool (i.e. either from the BOS/non-BOS unit population in each module). A single sample is
 664 denoted as $u_i^{b,n}$ where b is a Boolean variable indicating whether the ablated units are BOS units or non-
 665 BOS units, and $i = 1, 2, \dots, N_u$ represents the i^{th} unit sample. For each ablation sample $u_i^{b,n}$, the unit
 666 activity in the sample was set to zero. Mean-squared prediction error (MSE) was measured as the mean-
 667 squared difference between the predicted (\hat{A}_0) and actual frames (A_0), averaging over all pixels and time
 668 steps. Relative prediction error (RPE) of one video and one ablation sample was computed as

$$RPE(v_a^\alpha, u_i^{b,n}) \equiv RPE_{a,i}^{\alpha,b,n} = [MSE(v_a^\alpha, u_i^{b,n}) - MSE(v_a^\alpha, 0)] / MSE(v_a^\alpha, 0) \quad (5)$$

669 where $MSE(v_a^\alpha, 0)$ represents the MSE to the same video without ablation. We then computed the
 670 average RPE for a single video sample α and a given number of n ablated units:

$$RPE^{\alpha,b,n} = \langle RPE_{a,i}^{\alpha,b,n} \rangle_{a,i} \quad (6)$$

671 where $\langle \cdot \rangle_{a,i}$ represents the average across indices a and i . Dots and error bars in Figure 4 show the mean
 672 and SEM of $RPE^{\alpha,b,n}$ across different video samples α , with respect to the number of ablated units n , for
 673 the subsampled population (see previous Methods section: ‘Subsampling BOS and Non-BOS Units to
 674 Reduce Their Response Differences’). SI Figure 9 shows the result for the original population (without
 675 subsampling).

676

677 Statistical Analysis of the Ablation Experiment

678 We model the $RPE^{\alpha,b,n}$ as a linear model

$$RPE^{\alpha,b,n} = k^b n + \epsilon \quad (7)$$

679 where the intercept term is zero because the RPE is zero when no units are ablated. ϵ is an error term
 680 with a zero mean and a constant unknown variance, and k^b is the slope of a line that represents the
 681 average change in RPE if one additional unit is ablated ($b = bos$ for ablation of BOS units, $b = non -$
 682 bos for ablation of non-BOS units). We are interested in determining whether the slope k^{bos} is
 683 significantly different from $k^{non-bos}$. A bootstrap method is used as follows.

684 Observations are denoted as $D^b = \{RPE^{\alpha,b,n}\}$ where α and n indicate respectively video samples and
 685 number of ablated units. $N_s = 10,000$ bootstrap samples are generated by resampling D^b with
 686 replacement, denoted as $D^{b,s}$ where $s = 1, 2, \dots, N_s$. For each bootstrapped dataset, we used ordinary
 687 least squares linear regression to compute a slope $k^{b,s}$. 95% confidence interval of the slopes were
 688 estimated from the bootstrapped distribution (shown as error bands in Fig. 4 and SI Fig. 9). Subtracting
 689 the two slope sets, we got N_s slope differences denoted as $\Delta k^s = k^{bos,s} - k^{non-bos,s}$. The p-value (two-

690 tailed) was then estimated as $2 \times \min\{ Q(0, \{\Delta k^s\}), 1 - Q(0, \{\Delta k^s\}) \}$ where $Q(0, \{\Delta k^s\})$ is the
691 quantile of 0 in the set of slope differences $\{\Delta k^s\}$.
692

693 **Acknowledgments:**

694

695 **Funding:**

696 NIH grant R00EY031795 (TPF)

697 Incubator for Transdisciplinary Futures: Toward a Synergy Between Artificial Intelligence and
698 Neuroscience (RW).

699

700 **Author contributions:**

701 Conceptualization: ZY, RW, TPF

702 Methodology: ZY, RW, TPF

703 Investigation: ZY

704 Supervision: RW, TPF

705 Writing: ZY, RW, TPF

706 **Competing interests:** Authors declare that they have no competing interests

707

708 **Data and materials availability:**

709

710

711

712

713

714

715

References

- 716
717 1. Zhou, H., Friedman, H. S. & Von Der Heydt, R. Coding of border ownership in monkey visual
718 cortex. *J. Neurosci.* **20**, 6594–6611 (2000).
- 719 2. Hesse, X. J. K. & Tsao, D. Y. Consistency of Border-Ownership Cells across Artificial Stimuli ,
720 Natural Stimuli , and Stimuli with Ambiguous Contours. *J. Neurosci.* **36**, 11338–11349 (2016).
- 721 3. Franken, T. P. & Reynolds, J. H. Columnar processing of border ownership in primate visual
722 cortex. *Elife* **10**, 1–28 (2021).
- 723 4. Zhu, S., Oh, Y. J., Trepka, E., Chen, X. & Moore, T. Dependence of Contextual Modulation in
724 Macaque V1 on Interlaminar Signal Flow. *bioRxiv* 2024.04.18.590176 (2024)
725 doi:10.1101/2024.04.18.590176.
- 726 5. Williford, J. R. & Von Der Heydt, R. Figure-ground organization in visual cortex for natural
727 scenes. *eNeuro* **3**, 1–15 (2016).
- 728 6. Hu, B., von der Heydt, R. & Niebur, E. Figure-ground organization in natural scenes:
729 Performance of a recurrent neural model compared with neurons of area v2. *eNeuro* **6**, (2019).
- 730 7. Qiu, F. T. & Von Der Heydt, R. Figure and ground in the visual cortex: V2 combines stereoscopic
731 cues with Gestalt rules. *Neuron* **47**, 155–166 (2005).
- 732 8. Fang, F., Boyaci, H. & Kersten, D. Border ownership selectivity in human early visual cortex and
733 its modulation by attention. *J. Neurosci.* **29**, 460–465 (2009).
- 734 9. von der Heydt, R., Macuda, T. & Qiu, F. T. Border-ownership-dependent tilt aftereffect. *J. Opt.*
735 *Soc. Am. A* **22**, 2222 (2005).
- 736 10. Rideaux, R. & Harrison, W. J. Border ownership-dependent tilt aftereffect for shape defined by
737 binocular disparity and motion parallax. *J. Neurophysiol.* **121**, 1917–1923 (2019).
- 738 11. Victor, J. D. & Conte, M. M. Functional recursion of orientation cues in figure-ground separation.
739 *Vision Res.* **197**, 108047 (2022).
- 740 12. Kanwisher, N., Khosla, M. & Dobs, K. Using artificial neural networks to ask ‘why’ questions of
741 minds and brains. *Trends Neurosci.* **46**, 240–254 (2023).
- 742 13. Von der Heydt, R. Figure-ground organization and the emergence of proto-objects in the visual
743 cortex. *Front. Psychol.* **6**, 1–10 (2015).
- 744 14. LeCun, Y., Bottou, L., Bengio, Y. & Haffner, P. Gradient-based learning applied to document
745 recognition. *Proc. IEEE* **86**, 2278–2323 (1998).
- 746 15. Dedieu, A., Rikhye, R. V., Lázaro-Gredilla, M. & George, D. Learning attention-controllable
747 border-ownership for objectness inference and binding. *bioRxiv* 2020.12.31.424926 (2021).
- 748 16. Zhuang, C., Yan, S., Nayebi, A., Schrimpf, M., Frank, M. C., DiCarlo, J. J. & Yamins, D. L. K.
749 Unsupervised neural network models of the ventral visual stream. *Proc. Natl. Acad. Sci. U. S. A.*
750 **118**, (2021).
- 751 17. Lotter, W., Kreiman, G. & Cox, D. A neural network trained for prediction mimics diverse
752 features of biological neurons and perception. *Nat. Mach. Intell.* **2**, 210–219 (2020).
- 753 18. Eguchi, A. & Stringer, S. M. Neural network model develops border ownership representation
754 through visually guided learning. *Neurobiol. Learn. Mem.* **136**, 147–165 (2016).
- 755 19. Qiu, F. T., Sugihara, T. & Von Der Heydt, R. Figure-ground mechanisms provide structure for
756 selective attention. *Nat. Neurosci.* **10**, 1492–1499 (2007).
- 757 20. Oquab, M. *et al.* DINOv2: Learning Robust Visual Features without Supervision. 1–31 (2023).
- 758 21. Yamins, D. L. K., Hong, H., Cadieu, C. F., Solomon, E. A., Seibert, D. & DiCarlo, J. J.
759 Performance-optimized hierarchical models predict neural responses in higher visual cortex.
760 *Proc. Natl. Acad. Sci. U. S. A.* **111**, 8619–8624 (2014).
- 761 22. Richards, B. A. *et al.* A deep learning framework for neuroscience. *Nat. Neurosci.* **22**, 1761–1770

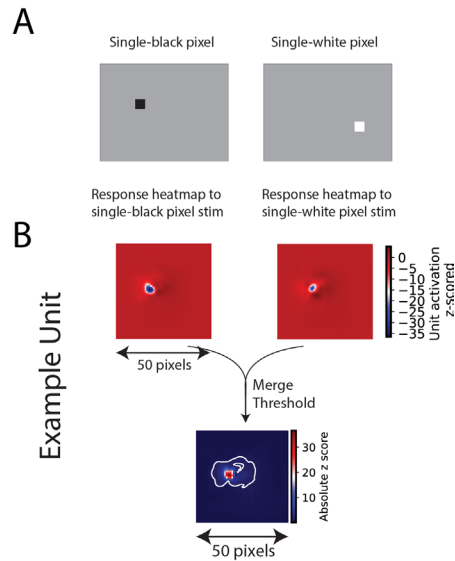
- 762 (2019).
- 763 23. Franken, T. P. & Reynolds, J. H. Grouping cells in primate visual cortex. *bioRxiv* 1–31 (2024)
- 764 doi:<https://doi.org/10.1101/2024.01.16.575953>.
- 765 24. O’Herron, P. & von der Heydt, R. Short-Term Memory for Figure-Ground Organization in the
- 766 Visual Cortex. *Neuron* **61**, 801–809 (2009).
- 767 25. O’Herron, P. & von der Heydt, R. Remapping of border ownership in the visual cortex. *J.*
- 768 *Neurosci.* **33**, 1964–1974 (2013).
- 769 26. Lotter, W., Kreiman, G. & Cox, D. Deep predictive coding networks for video prediction and
- 770 unsupervised learning. *5th Int. Conf. Learn. Represent. ICLR 2017 - Conf. Track Proc.* 1–18
- 771 (2017).
- 772 27. Rane, R. P., Szügyi, E., Saxena, V., Ofner, A. & Stober, S. PredNet and predictive coding: A
- 773 critical review. *ICMR 2020 - Proc. 2020 Int. Conf. Multimed. Retr.* 233–241 (2020)
- 774 doi:10.1145/3372278.3390694.
- 775 28. Geiger, A., Lenz, P., Stiller, C. & Urtasun, R. Vision meets robotics: The KITTI dataset. *Int. J.*
- 776 *Rob. Res.* 1–6 (2013).
- 777 29. Pospisil, D. A., Pasupathy, A. & Bair, W. ‘Artphysiology’ reveals V4-like shape tuning in a deep
- 778 network trained for image classification. *Elife* **7**, 1–30 (2018).
- 779 30. von der Heydt, R. Visual cortical processing—From image to object representation. *Front.*
- 780 *Comput. Sci.* **5**, (2023).
- 781 31. Zhang, N. R. & Von Der Heydt, R. Analysis of the context integration mechanisms underlying
- 782 figure-ground organization in the visual cortex. *J. Neurosci.* **30**, 6482–6496 (2010).
- 783 32. Nakayama, K., He, Z. J. & Shimojo, S. Visual surface representation: A critical link between
- 784 lower-level and higher-level vision. in *In Invitation to Cognitive Science* 1–70 (The MIT Press,
- 785 1995).
- 786 33. Rao, R. P. N. & Ballard, D. H. Predictive coding in the visual cortex: a functional interpretation
- 787 of some extra-classical receptive-field effects. *Nat. Neurosci.* **2**, 79–87 (1999).
- 788 34. Friston, K. The free-energy principle: A unified brain theory? *Nat. Rev. Neurosci.* **11**, 127–138
- 789 (2010).
- 790 35. Shipp, S. Neural elements for predictive coding. *Front. Psychol.* **7**, 1–21 (2016).
- 791 36. Bastos, A. M., Usrey, W. M., Adams, R. A., Mangun, G. R., Fries, P. & Friston, K. J. Canonical
- 792 Microcircuits for Predictive Coding. *Neuron* **76**, 695–711 (2012).
- 793 37. Hosseini, M. & Maida, A. Hierarchical Predictive Coding Models in a Deep-Learning
- 794 Framework. *arXiv* (2020) doi:10.48550/arXiv.2005.03230.
- 795 38. Jehee, J. F. M., Lamme, V. A. F. & Roelfsema, P. R. Boundary assignment in a recurrent network
- 796 architecture. *Vision Res.* **47**, 1153–1165 (2007).
- 797 39. Mehrani, P. & Tsotsos, J. K. Early recurrence enables figure border ownership. *Vision Res.* **186**,
- 798 23–33 (2021).
- 799 40. Zhu, S. D., Zhang, L. A. & Von Der Heydt, R. Searching for object pointers in the visual cortex.
- 800 *J. Neurophysiol.* **123**, 1979–1994 (2020).
- 801 41. Craft, E., Schütze, H., Niebur, E. & Von Der Heydt, R. A neural model of figure-ground
- 802 organization. *J. Neurophysiol.* **97**, 4310–4326 (2007).
- 803 42. Martin, A. B. & von der Heydt, R. Spike synchrony reveals emergence of proto-objects in visual
- 804 cortex. *J. Neurosci.* **35**, 6860–6870 (2015).
- 805 43. Zhaoping, L. Border ownership from intracortical interactions in visual area V2. *Neuron* **47**, 143–
- 806 153 (2005).
- 807 44. Sakai, K., Nishimura, H., Shimizu, R. & Kondo, K. Consistent and robust determination of border

- 808 ownership based on asymmetric surrounding contrast. *Neural Networks* **33**, 257–274 (2012).
- 809 45. Sakai, K. & Nishimura, H. Surrounding suppression and facilitation in the determination of
810 border ownership. *J. Cogn. Neurosci.* **18**, 562–579 (2006).
- 811 46. Supèr, H., Romeo, A. & Keil, M. Feed-forward segmentation of figure-ground and assignment of
812 border-ownership. *PLoS One* **5**, (2010).
- 813 47. O’Herron, P. & von der Heydt, R. Representation of object continuity in the visual cortex. *J. Vis.*
814 **11**, 1–9 (2011).
- 815 48. Heider, B., Spillmann, L. & Peterhans, E. Stereoscopic illusory contours - Cortical neuron
816 responses and human perception. *J. Cogn. Neurosci.* **14**, 1018–1029 (2002).
- 817 49. Yang, G. R. & Wang, X. J. Artificial Neural Networks for Neuroscientists: A Primer. *Neuron*
818 **107**, 1048–1070 (2020).
- 819 50. Taylor, J. M. & Kriegerkorte, N. Extracting and visualizing hidden activations and computational
820 graphs of PyTorch models with TorchLens. *Sci. Rep.* **13**, 1–13 (2023).
- 821 51. Haeusler, S. & Maass, W. A statistical analysis of information-processing properties of lamina-
822 specific cortical microcircuit models. *Cereb. Cortex* **17**, 149–162 (2007).
- 823 52. Douglas, R. J. & Martin, K. A. A functional microcircuit for cat visual cortex. *J. Physiol.* **440**,
824 735–769 (1991).
- 825 53. Oprea, S., Martinez-Gonzalez, P., Garcia-Garcia, A., Castro-Vargas, J. A., Orts-Escolano, S.,
826 Garcia-Rodriguez, J. & Argyros, A. A Review on Deep Learning Techniques for Video
827 Prediction. *IEEE Trans. Pattern Anal. Mach. Intell.* 1–26 (2020)
828 doi:10.1109/TPAMI.2020.3045007.
- 829 54. Bear, D. M., Feigelis, K., Chen, H., Lee, W., Venkatesh, R., Kotar, K., Durango, A. & Yamins,
830 D. L. K. Unifying (Machine) Vision via Counterfactual World Modeling. *arXiv* 1–22 (2023)
831 doi:10.48550/arXiv.2306.01828.
- 832 55. Wu, Y., Gao, R., Park, J. & Chen, Q. Future video synthesis with object motion prediction. *Proc.*
833 *IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.* 5538–5547 (2020)
834 doi:10.1109/CVPR42600.2020.00558.
- 835 56. Liu, Z., Yeh, R. A., Tang, X., Liu, Y. & Agarwala, A. Video Frame Synthesis Using Deep Voxel
836 Flow. *Proc. IEEE Int. Conf. Comput. Vis.* **2017-October**, 4473–4481 (2017).
- 837 57. Finn, C., Goodfellow, I. & Levine, S. Unsupervised learning for physical interaction through
838 video prediction. *Adv. Neural Inf. Process. Syst.* 64–72 (2016).
- 839 58. Tulyakov, S., Liu, M. Y., Yang, X. & Kautz, J. MoCoGAN: Decomposing Motion and Content
840 for Video Generation. *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.* 1526–1535
841 (2018) doi:10.1109/CVPR.2018.00165.
- 842 59. Lee, S., Kim, H. G., Choi, D. H., Kim, H. Il & Ro, Y. M. Video Prediction Recalling Long-term
843 Motion Context via Memory Alignment Learning. *Proc. IEEE Comput. Soc. Conf. Comput. Vis.*
844 *Pattern Recognit.* 3053–3062 (2021) doi:10.1109/CVPR46437.2021.00307.
- 845 60. Zhong, Y., Liang, L., Zharkov, I. & Neumann, U. MMVP: Motion-Matrix-based Video
846 Prediction. *Proc. IEEE Int. Conf. Comput. Vis.* 4250–4260 (2023)
847 doi:10.1109/ICCV51070.2023.00394.
- 848 61. Denton, E. & Birodkar, V. Unsupervised learning of disentangled representations from video.
849 *Adv. Neural Inf. Process. Syst.* **2017-December**, 4415–4424 (2017).
- 850 62. Geirhos, R., Jacobsen, J. H., Michaelis, C., Zemel, R., Brendel, W., Bethge, M. & Wichmann, F.
851 A. Shortcut learning in deep neural networks. *Nat. Mach. Intell.* **2**, 665–673 (2020).
- 852 63. Geirhos, R., Michaelis, C., Wichmann, F. A., Rubisch, P., Bethge, M. & Brendel, W. Imagenet-
853 trained CNNs are biased towards texture; increasing shape bias improves accuracy and

- 854 robustness. *7th Int. Conf. Learn. Represent. ICLR 2019* 1–22 (2019).
- 855 64. Beery, S., Van Horn, G. & Perona, P. Recognition in Terra Incognita. *Lect. Notes Comput. Sci.*
856 *(including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)* **11220 LNCS**, 472–489
857 (2018).
- 858 65. Baker, N., Lu, H., Erlikhman, G. & Kellman, P. J. *Deep convolutional networks do not classify*
859 *based on global object shape. PLoS Computational Biology* vol. 14 (2018).
- 860 66. Biederman, I. & Ju, G. Surface versus edge-based determinants of visual recognition. *Cogn.*
861 *Psychol.* **20**, 38–64 (1988).
- 862 67. Luongo, F. J., Liu, L., Ho, C. L. A., Hesse, J. K., Wekselblatt, J. B., Lanfranchi, F. F., Huber, D.
863 & Tsao, D. Y. Mice and primates use distinct strategies for visual segmentation. *Elife* **12**, 1–30
864 (2023).
- 865 68. Wilson, E. B. Probable Inference, the Law of Succession, and Statistical Inference. *J. Am. Stat.*
866 *Assoc.* **22**, 209–212 (1927).
- 867

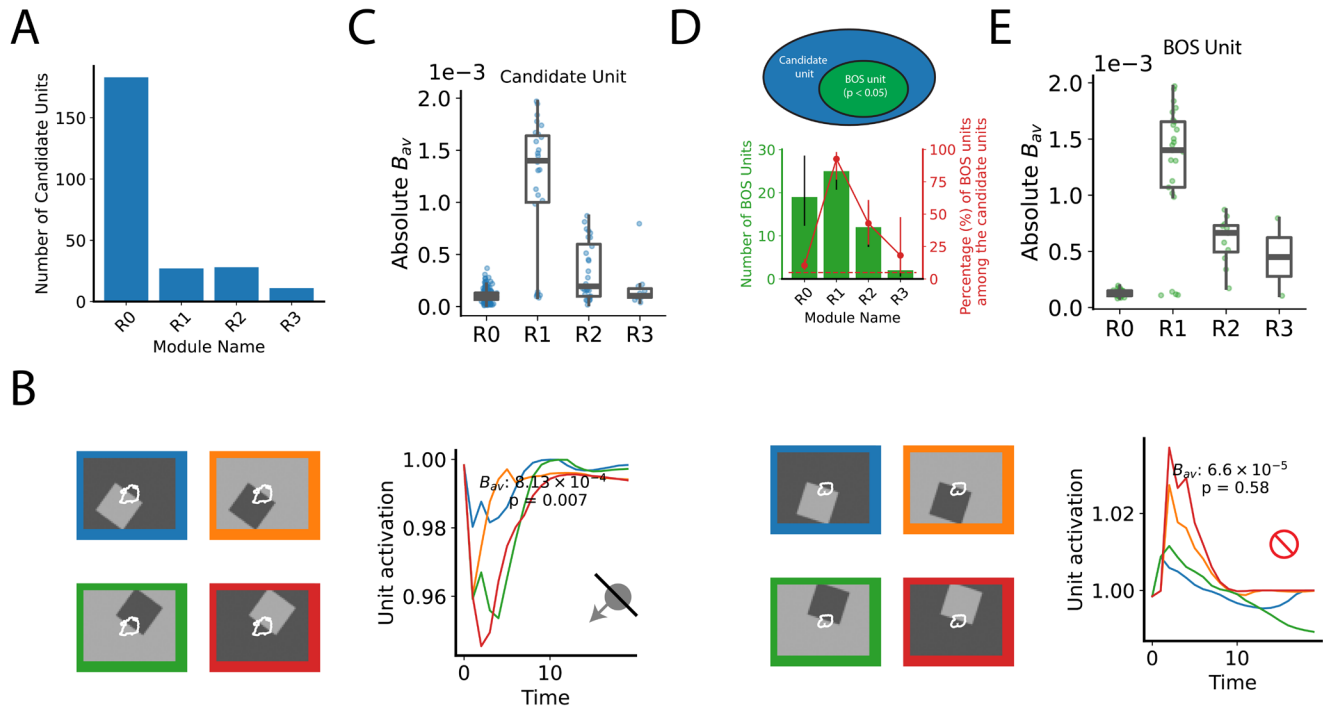
868
869
870
871

Supplementary Materials



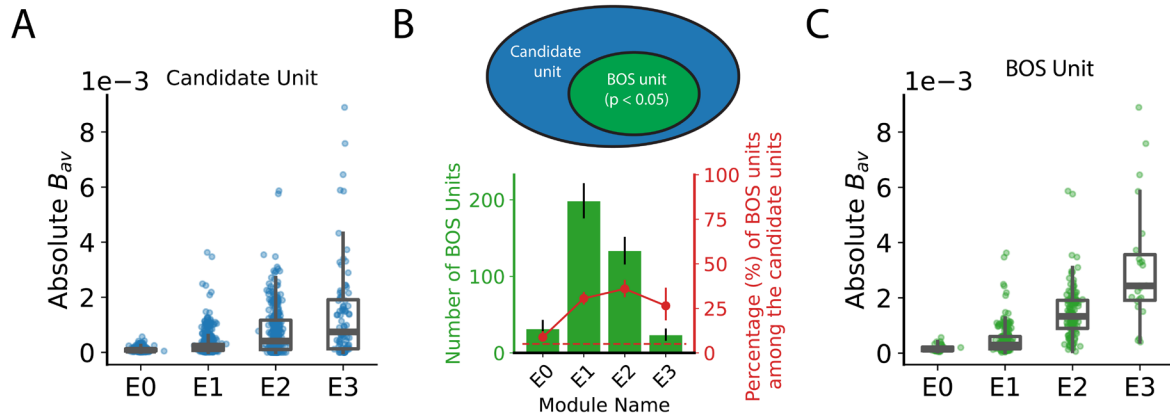
872
873
874
875
876
877
878
879
880
881

SI Figure 1. Illustration of the method to measure the cRF of PredNet units. (A) A sparse noise scene is a gray scene with only one black or white pixel, at a random position. These scenes were used as input to PredNet over four time steps. (B) cRF for an example unit. The unit's responses to the sparse noise scenes were collected and normalized (z-scored) into two heat maps, one for black pixel noise and the other for white. Each value in the black or white heatmap corresponds to the unit's normalized response to a black or white pixel at the same entry position. The two heatmaps (for one unit) were merged into one heatmap by taking the maximum absolute values for each entry. Positions with an absolute value of the z-score greater than 1 were defined as the cRF (indicated by white contours).



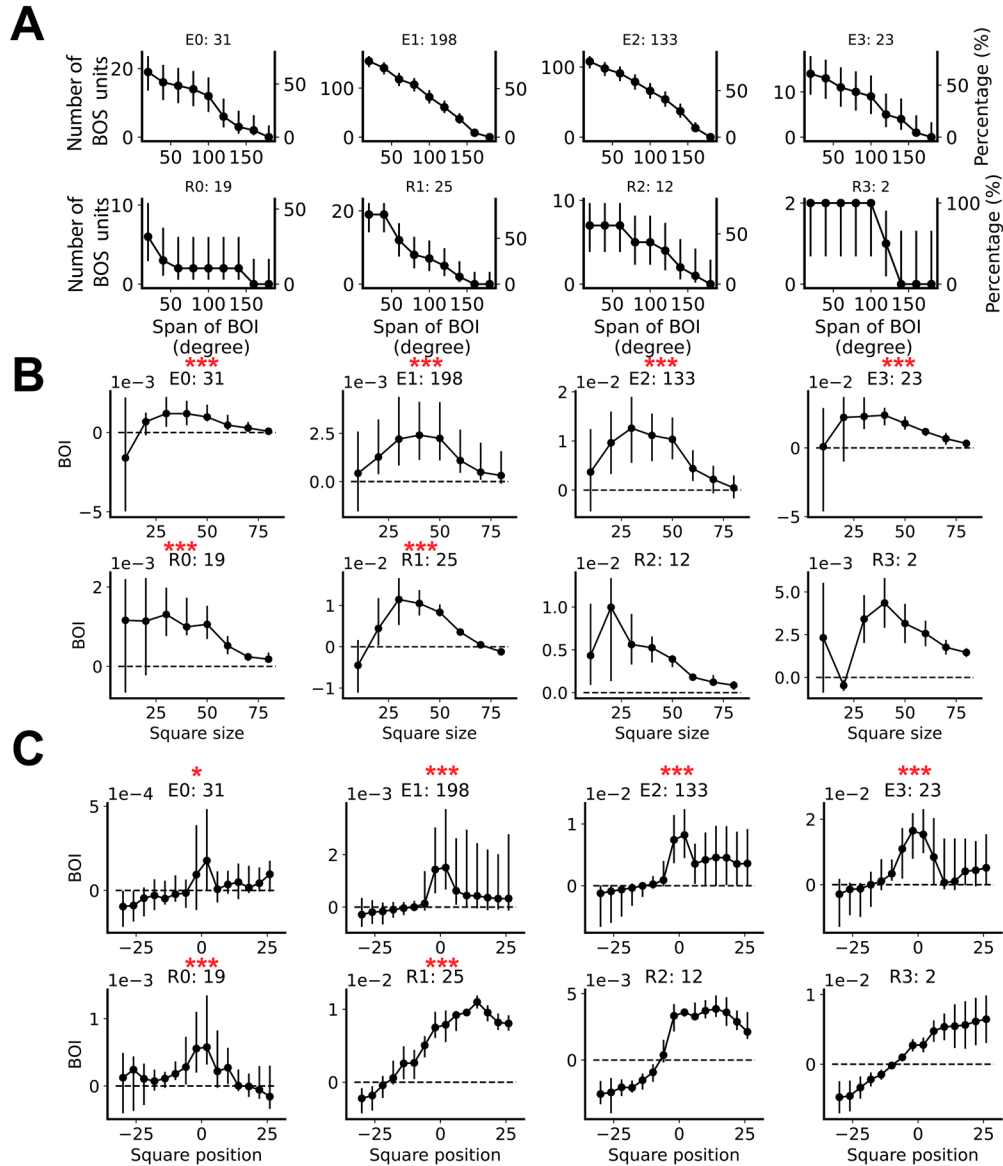
882
883
884
885
886
887
888
889
890
891
892

SI Figure 2. BOS units emerge in PredNet's R modules. (A) The number of candidate units in R modules across different layers. (B) Responses of two example units (module R_2), with white contours indicating the cRF (similar to Fig. 1D). B_{av} measures the unit's response different to different BOSs across different square orientations (see Methods). P value (two tailed) was computed by comparing B_{av} to that after shuffling stimulus labels (permutation test, see Methods). Arrow in the middle-left panel indicates the preferred side of BOS for the example candidate unit. (C) The B_{av} distribution of the candidate units in different R modules. Each dot is one candidate unit. (D) Among the candidate units, units with p-value smaller than 0.05 are defined as BOS units. Error bars indicate 95% confidence intervals. Horizontal dashed line indicates chance level of 5%. (E) The B_{av} distribution of BOS units in different R modules. Each dot is one BOS unit.



893
894
895

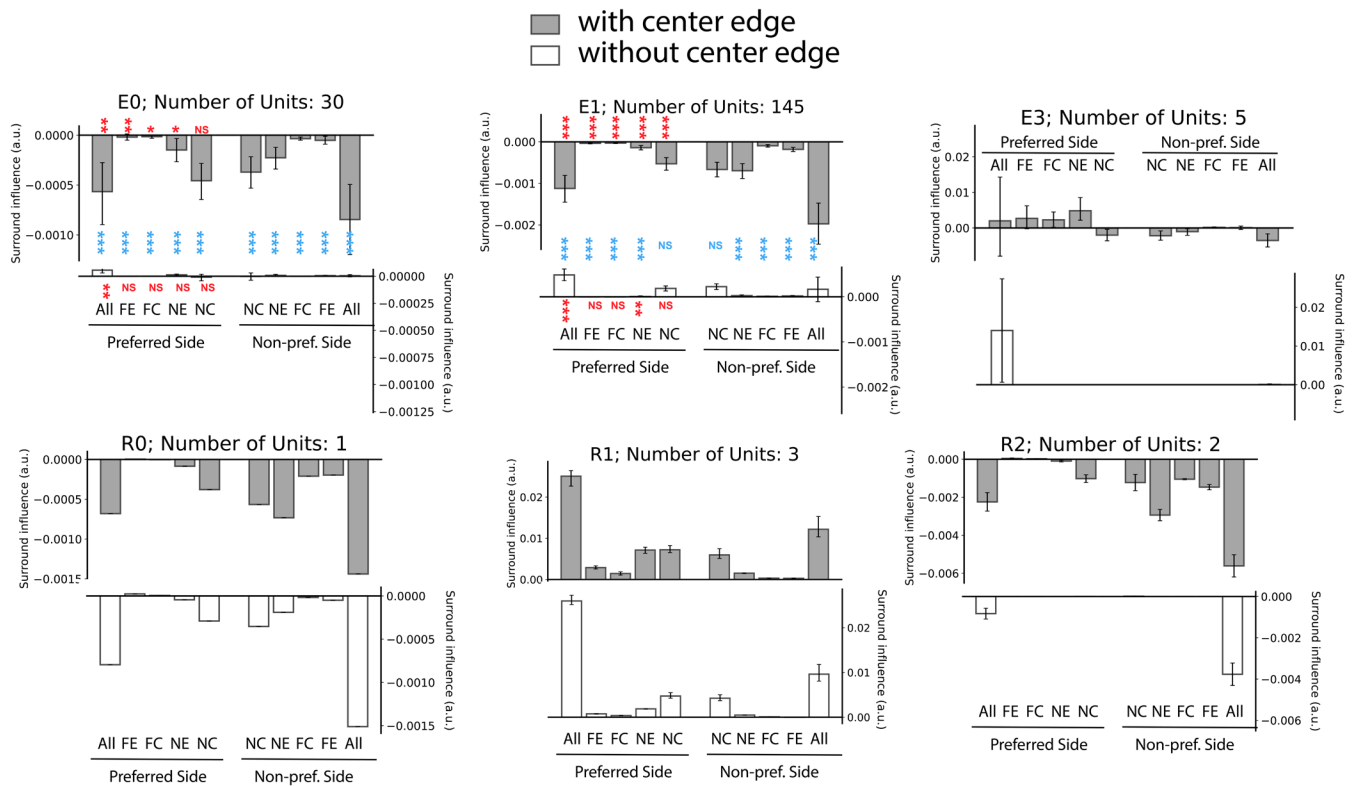
SI Figure 3. B_{av} distribution of units in E modules. Similar as SI Fig. 2C-E, for E modules.



896
897
898
899
900
901
902
903
904

SI Figure 4. BOS signals are robust to different stimulus parameters. (A) Similar to Figure 2B, for other modules. (B, C) BOS across different square sizes and positions. The dots and error bars represent the median, first and third quartiles across all units in a module. The number after the module name in the panel titles denotes the total number of BOS units included per module. Red symbols indicate whether the averaged BOS across conditions (square sizes or positions) are statistically significantly larger than zero, ***: $p < 0.001$; *: $p < 0.05$; bootstrapping test (see Methods). Statistical significance was only evaluated in modules with more than 15 BOS units.

905



906

907

908

909

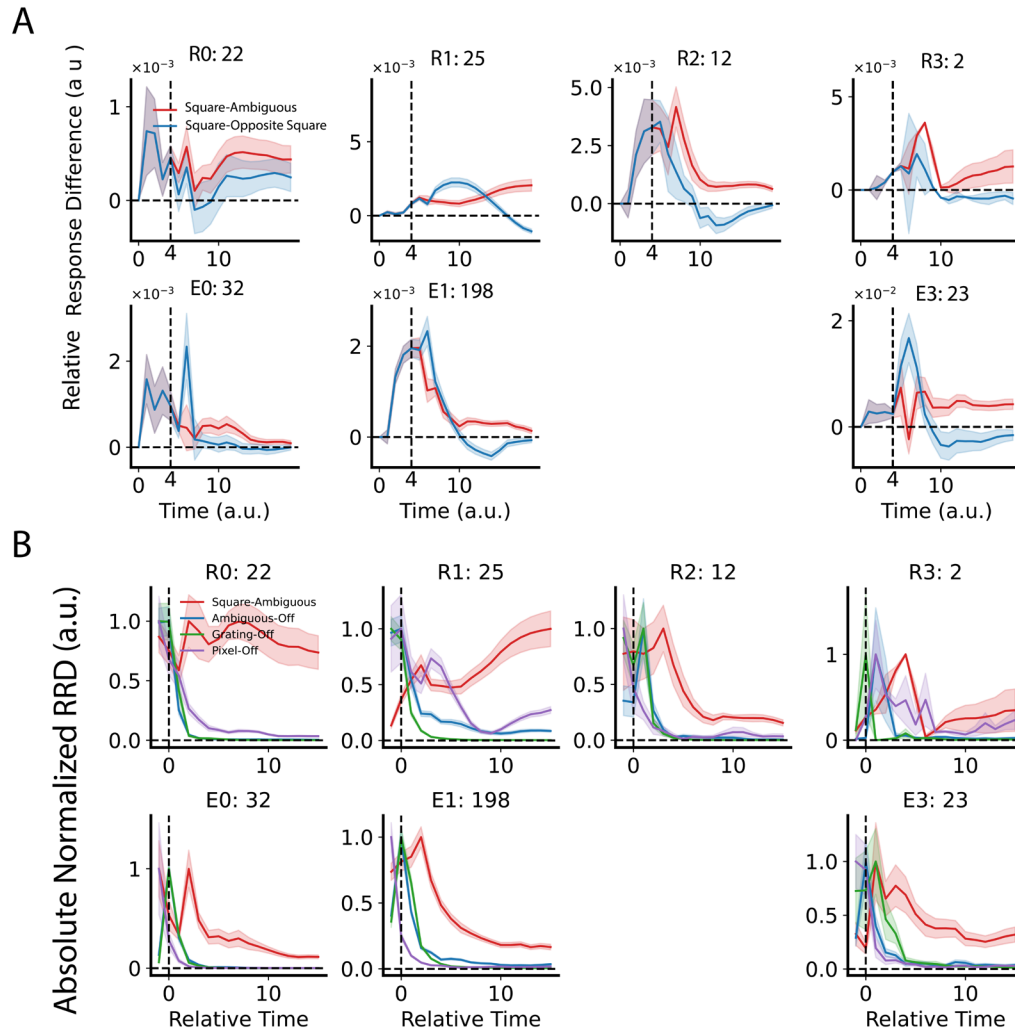
910

911

912

913

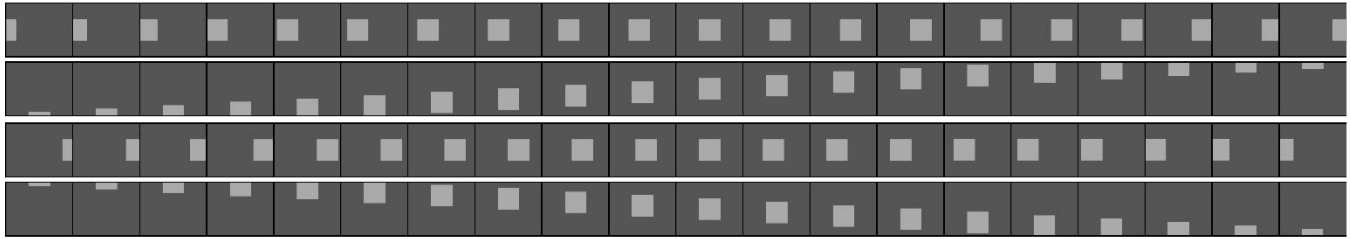
SI Figure 5. BOS units' responses to square fragments on the preferred side of BOS are generally larger than on the non-preferred side of BOS. Similar to Fig. 2F, for BOS units in different PredNet modules. Red text indicates whether the surround influence for a particular condition is significantly larger on the preferred side than on the non-preferred side. Blue text indicates whether the absolute value of surround influence of with-CE is significantly larger than without-CE case. Wilcoxon signed-rank test. ***: $p < 0.001$; **: $p < 0.01$; *: $p < 0.05$; NS: no significance.



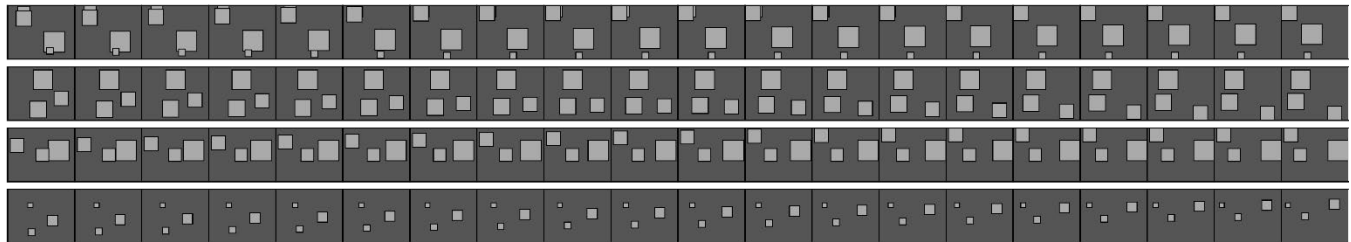
914
915
916
917

SI Figure 6. Persistent BOS signals in different modules. (A) Similar as Fig. 3B, for other modules. The number of BOS units in each module is indicated in the title. (B) Similar as Fig. 3C, for other modules.

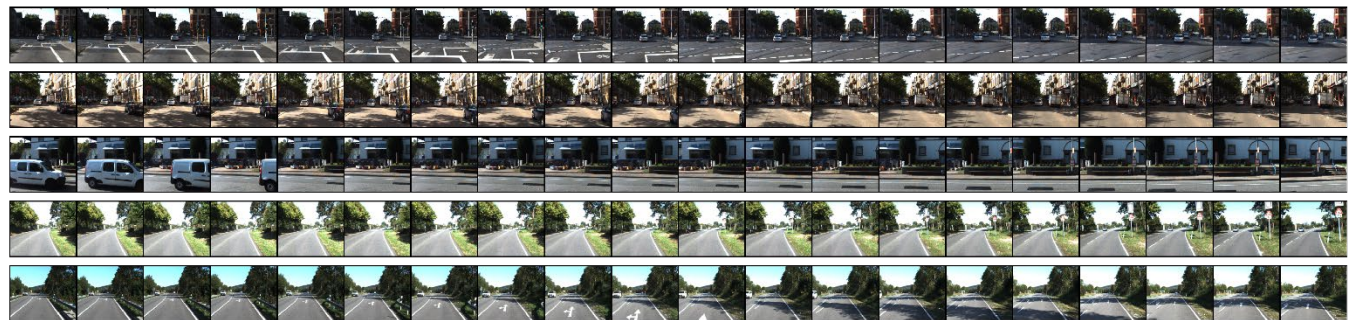
Translating Video Examples



Random Squares Video Examples



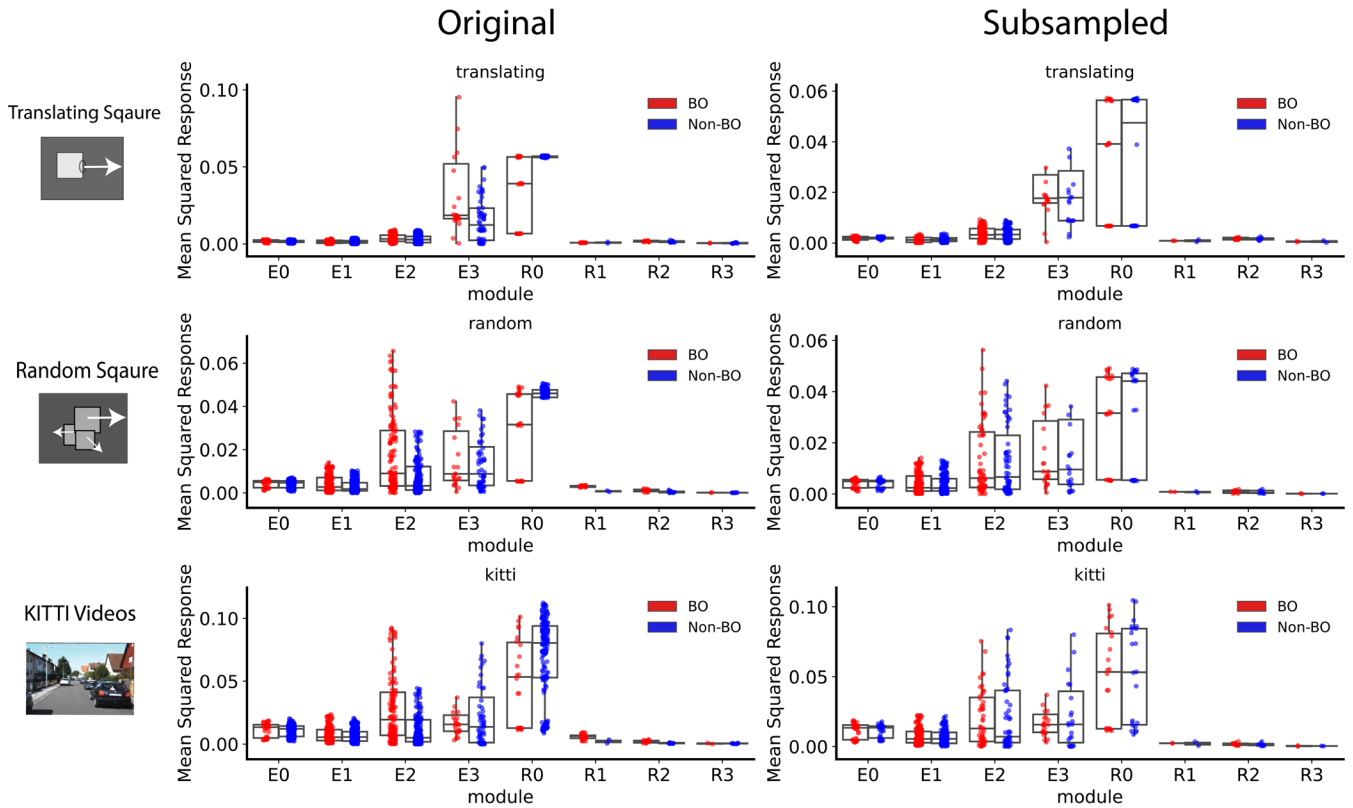
KITTI Video Examples



918
919
920
921
922

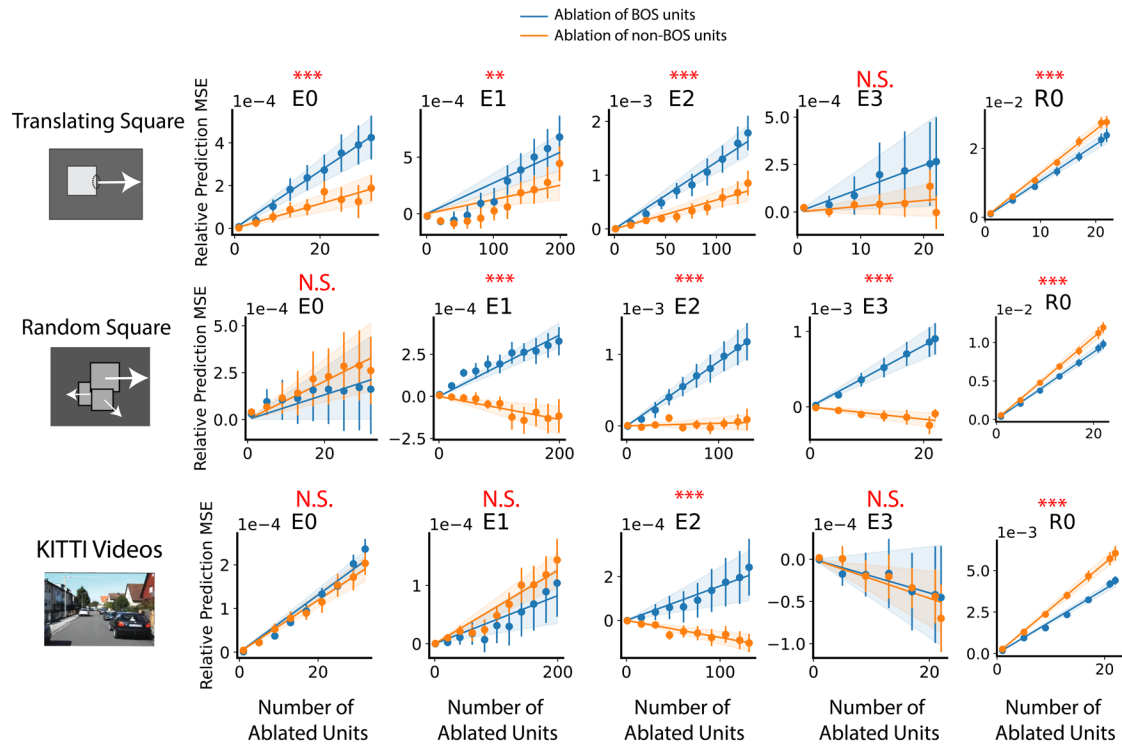
SI Figure 7. Three video types used for the ablation experiment. Figure shows example videos from each video type. Each row shows a different unique video for each of the three types. Video length is 20 frames, shown during 20 time steps.

923
924



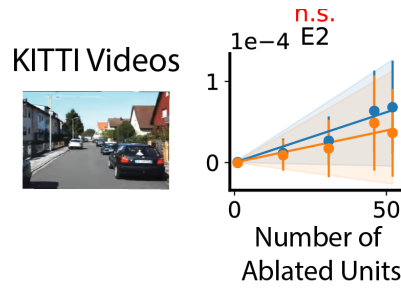
925
926
927
928
929
930
931
932

SI Figure 8. Activity in subsampled BOS and non-BOS unit populations and original populations. For each unit, mean squared response is the square of the averaged response, averaging cross time and videos. Each dot is one unit's mean squared response. Boxes indicate the interquartile range between the first and third quartiles with central mark inside each box indicating the median. Whiskers extend to the lowest and highest values within 1.5 times the interquartile range. Outlier units not shown for better visualization (but included in the metrics indicated by the boxplots).



933
934
935
936

SI Figure 9. The effect of ablating the original BOS/non-BOS units, without subsampling. Similar as Figure 4 but using original unit population (no subsampling).



937
938
939
940

SI Figure 10. The effect of ablating the subsampled E_2 BOS/non-BOS units on KITTI video prediction. Similar as Figure 4, for KITTI videos.