

---

# Unifying Gestalt Principles Through Inference-Time Prior Integration

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Gestalt principles such as closure, similarity, continuation, and figure–ground  
2 segregation have often been characterized as distinct rules guiding perceptual  
3 organization. We demonstrate that these diverse phenomena emerge from a single  
4 computational mechanism: inference-time prior integration. Our algorithm, Prior-  
5 Guided Drift Diffusion (PGDD), repurposes the same feedback pathways used  
6 during backpropagation to refine neural activations during inference, enabling  
7 networks to integrate learned statistical regularities with sensory input. Applied to  
8 pre-trained networks, PGDD reproduces illusory contours, perceptual grouping, and  
9 figure-ground segregation without additional training or architectural modifications.  
10 These effects depend on appropriate learned priors, for example, networks trained  
11 on object-centric datasets generate Kanizsa illusions, while those trained only on  
12 faces or scenes fail to do so. Our results suggest that Gestalt principles are not  
13 hardwired perceptual rules but emergent consequences of how neural networks  
14 can dynamically combine learned statistical knowledge with incoming sensory  
15 evidence. This computational framework bridges artificial and biological vision by  
16 showing how inference-time optimization can account for fundamental aspects of  
17 human perceptual organization.

## 18 1 Introduction

19 Gestalt psychology identified a set of robust phenomena of perceptual organization: closure, similarity,  
20 continuation, and figure–ground segregation [1]. These effects have been central to vision science  
21 for more than a century. They demonstrate that perception is not a direct reflection of the stimulus,  
22 but instead reflects internal constraints that group elements into coherent wholes. Several theoretical  
23 frameworks have attempted to explain these effects. Structural information theory formalized Gestalt  
24 laws through preferences for simpler encodings [2]. Bayesian approaches treated grouping as  
25 probabilistic inference over common causes [3]. Predictive coding models proposed that feedback  
26 carries expectations to generate illusory contours and modulate grouped responses [4, 5]. Despite  
27 these efforts, no existing account provides a single mechanistic explanation that unifies all classical  
28 Gestalt principles within the same computational process that also underlies typical visual perception.

29 Neural networks trained for object recognition show alignment with human behavioral and neural  
30 activity patterns on natural images [6, 7]. However, they still fall short on many perceptual phenomena  
31 tied to Gestalt principles. Prior studies have found that deep CNNs often fail to perceive illusory  
32 contours and do not reliably discriminate based on global shape cues, instead relying more heavily  
33 on local texture or outlines [8–10]. Some Gestalt-like effects have been reported in convolutional  
34 neural networks. For example, Kim et al. [11] showed that closure emerges spontaneously in  
35 networks trained on natural images, and Biscione & Bowers [12] reported mixed evidence for  
36 grouping. However, these effects do not follow the rules observed in biological systems. In the brain,

37 Gestalt phenomena such as illusory contours and figure–ground segregation are characterized by  
 38 delayed neural responses, and causal dependence on intact feedback pathways [1, 13–15]. While  
 39 representations in standard feedforward CNNs can sometimes be interpreted as Gestalt-like, they do  
 40 not exhibit the temporal delays or feedback dependence observed in biological systems, suggesting  
 41 that they may rely on a different underlying mechanism.

42 Inspired by the biological evidence that Gestalt phenomena depend on cortical feedback (Fig. 1  
 43 A), we hypothesized that the same feedback signals used during learning could also be used during  
 44 inference. In this work, we propose that the feedback pathways used to propagate error gradients  
 45 during backpropagation may also serve as a channel for accessing the implicit priors a classifier  
 46 acquires during training. We devised an algorithm, *Prior-Guided Drift Diffusion* (PGDD), that  
 47 reuses these learning feedback signals to update activations at inference time, without any additional  
 48 training or new feedback connections. PGDD applied to both convolutional and transformer pre-  
 49 trained networks reproduces classic Gestalt demonstrations—including Kanizsa figures, similarity  
 50 grouping, contour continuation, and figure–ground segregation. Each of these specific stimuli has  
 51 been the subject of decades of psychological and neurophysiological experiments [1, 13, 14, 16–20],  
 52 providing a direct link between our computational results and established experimental findings. This  
 53 demonstrates that all of these phenomena can be unified under a single mechanism: the integration of  
 54 learned priors with incoming sensory input.

## 55 2 Prior-Guided Drift Diffusion Algorithm

56 Our central hypothesis is that the same feedback signals used for learning can be repurposed during  
 57 inference to integrate priors with incoming sensory input. To test this idea in generic neural networks,  
 58 we developed an inference procedure that makes this mechanism explicit. We call this procedure *Prior-  
 59 Guided Drift Diffusion* (PGDD). PGDD refines activations iteratively by drifting away from noisy or  
 60 incomplete representations and diffusing toward states consistent with the network’s learned data  
 61 distribution (Fig. 1 B). We implemented PGDD as a feedback-mediated inference-time algorithm to  
 62 update activations. Given a network  $f_\theta$  trained on natural images under a robust-to-noise classification  
 63 objective, PGDD iteratively refines activations to align with learned priors:

$$\mathcal{L}_{\text{PGDD}} = \|r(x_t) - r(x_{\text{noisy}})\|^2 \quad (1)$$

64 where  $r(\cdot)$  denotes representations at a chosen layer,  $x_{\text{noisy}} = x + \epsilon$  with  $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$ , and the  
 65 constraint  $\|x_{t+1} - x_{\text{sensory}}\| \leq \epsilon_f$  prevents runaway updates. The update rule is:

$$x_{t+1} = x_t + \eta \nabla_{x_t} \mathcal{L}_{\text{PGDD}} + \xi_t, \quad \xi_t \sim \mathcal{N}(0, I). \quad (2)$$

66 This procedure encourages activations to move toward regions of representation space that are more  
 67 consistent with the network’s training distribution while maintaining fidelity to sensory input. A  
 68 theoretical analysis, provided in the Supplementary (Section A.1), shows how the PGDD update  
 69 connects formally to gradients of the data distribution, providing a principled justification for its  
 70 ability to recover Gestalt-like structure.

## 71 3 Results

72 We presented the exact same stimuli used in decades of psychological and neurophysiological  
 73 experiments on Gestalt perception. Unless otherwise noted, all results were obtained with a robustly  
 74 trained ResNet-50 on ImageNet [21]; in the Supp.A.4, we further show that a probabilistic variant of  
 75 PGDD applied to standard (non-robust) networks produces qualitatively similar outcomes. **Gestalt  
 76 Grouping: Similarity and Continuation** PGDD implemented classical Gestalt grouping principles  
 77 tested in monkeys [14] through the same mechanism (Fig. 1 C,D). For similarity grouping, elements  
 78 sharing color or shape characteristics were progressively linked through iterative feedback updates.  
 79 Good continuation was demonstrated using interrupted line segments. PGDD completed contours  
 80 across gaps, extending fragmented visual information into smooth, continuous patterns.

81 **Figure-ground segregation principle** In primate V1, neurons show enhanced responses to figure  
 82 regions compared to uniform backgrounds even when local features are identical [13]. PGDD  
 83 reproduced this texture-defined figure-ground effect by progressively amplifying boundaries while  
 84 suppressing background patterns, matching both the spatial selectivity and temporal dynamics  
 85 observed in cortical recordings. PGDD applied to face-vase illusion, often framed as an instance of

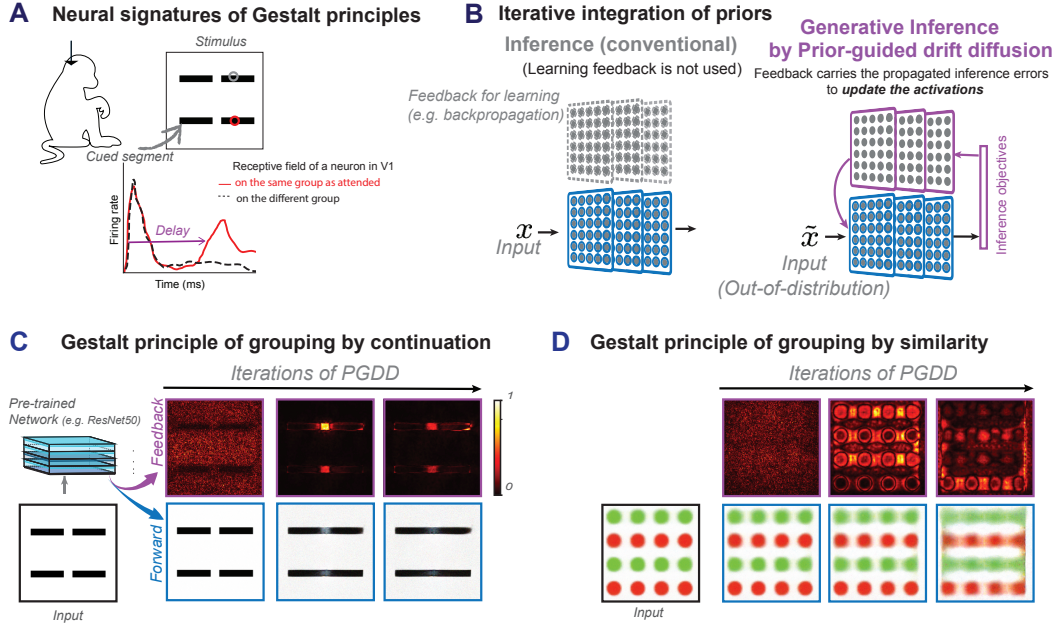


Figure 1: PGDD unifies Gestalt principles through feedback-driven prior integration. (A) Neural signatures of Gestalt principles in biological vision show characteristic delayed responses when neurons' receptive fields fall on grouped versus ungrouped elements, indicating feedback-mediated processing [1, 17], see also [13, 14, 16, 19, 20]. (B) PGDD repurposes learning feedback pathways during inference. While conventional inference leaves learning feedback connections unused, PGDD reactivates these same pathways (e.g., backpropagation) to iteratively update activations using inference objectives, enabling integration of learned priors with out-of-distribution inputs (Supp.A.2). (C) Gestalt principle of grouping by continuation: PGDD applied to the same stimulus in [1] progressively completes segments across iterations, with feedback activation patterns (top row) showing enhanced connectivity while forward activations (bottom row) display completed continuous lines. (D) Gestalt principle of grouping by similarity: PGDD groups elements sharing visual features (color, shape) through iterative feedback updates (same stimulus used in [1]).

86 figure-ground segregation, revealed how different learned priors shape perceptual interpretation (Fig.  
 87 2) B. Networks trained for object recognition generated vase-like features through PGDD, while those  
 88 trained for face recognition produced face-like patterns. Control experiments with simple rectangular  
 89 shapes showed that both appropriate sensory input and relevant priors are required for structured  
 90 interpretation.

91 **Closure principle: Kanizsa Figures** When presented with Kanizsa triangle inducers, PGDD gener-  
 92 ated clear illusory contours completing the missing edges (Figure 2). Initial feedforward processing  
 93 showed no evidence of triangle completion, but iterative prior accumulation progressively structured  
 94 activations to represent the completed figure. Control experiments confirmed the prior-dependence of  
 95 this effect. Networks trained on faces or places failed to generate Kanizsa squares, while those trained  
 96 on object-centric datasets succeeded, highlighting the role of learned statistical regularities. Across  
 97 all tested Gestalt principles, PGDD operated through the same computational process: iterative  
 98 accumulation of learned statistical regularities through feedback-mediated inference. The specific  
 99 manifestation depended on the network's training history and the choice of layer for implementing  
 100 the PGDD parameters (Supp. A.3), but the underlying mechanism remained the same.

101 **Limitations** A key constraint in our current implementation stems from using pretrained neural  
 102 networks (both robust and standard), which restricts our activation updates to the input layer. In the  
 103 brain, sensory input remains unchanged while early cortical areas show delayed, feedback-driven  
 104 responses. Our approach of updating input activations serves as a computational proxy for this  
 105 biological process. We predict that training neural networks where perturbations are applied to initial  
 106 layer (rather than inputs) would yield similar Gestalt phenomena while more closely matching the  
 107 biological implementation. This architectural modification would allow internal representations to

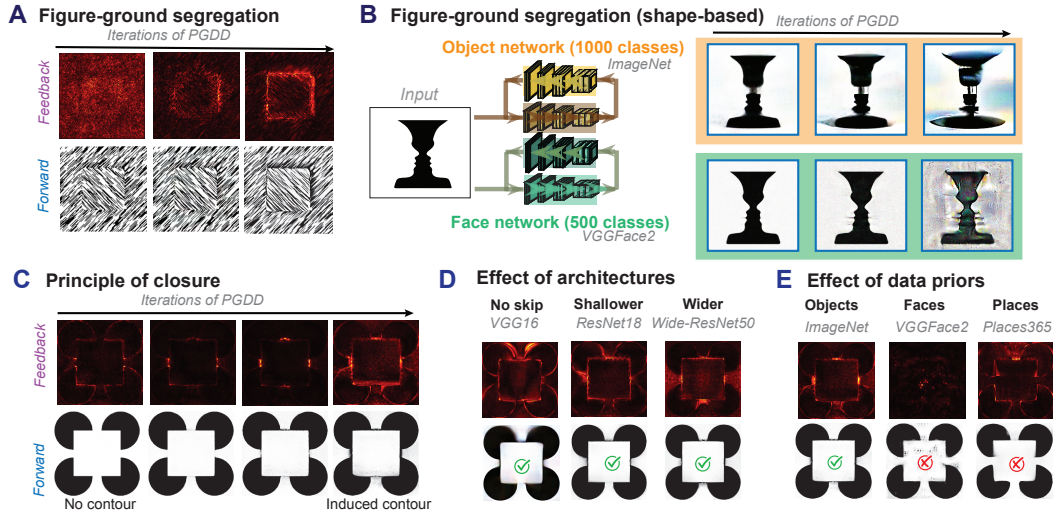


Figure 2: PGDD accounts for instances of Gestalt principles of figure-ground segregation and closure. (A) Figure-ground segregation: PGDD applied to texture-defined stimuli (the same stimulus in [13]) progressively enhances figure boundaries through iterative feedback updates, consistent with delayed V1 responses observed in primate studies. (B) Prior-dependent interpretation demonstrates the role of learned statistical regularities in ambiguous perception. When presented with the Rubin face-vase [18] stimulus, networks with different training regimens produce distinct interpretations: object networks (ImageNet, top) generate vase-like interpretations while face networks (VGGFace2, bottom) produce face-like patterns through the same algorithm (PGDD) accounting for perception-dependent activity in human’s early visual areas [18]. (C) Closure through prior integration: Kanizsa square inducers processed by PGDD show progressive emergence of illusory contours. Feedback activations (top) reveal the iterative strengthening of missing edges, while forward activations (bottom) display the completed square structure emerging over iterations, matching the delayed responses to illusory contours in early visual cortex [15]. (D) Architectural robustness: PGDD produces consistent closure effects across different network architectures (VGG16, ResNet18, Wide-ResNet50), demonstrating that the mechanism operates independently of specific architectural features. (E) Prior-specificity: Closure effects depend critically on appropriate learned priors. Networks trained on object-rich datasets (ImageNet) successfully complete Kanizsa figures, while networks trained exclusively on faces (VGGFace2) or scenes (Places365) fail to generate illusory contours, confirming that statistical regularities from training data drive perceptual completion.

108 be refined through feedback while maintaining fixed sensory input, better aligning with the layered  
 109 processing observed in visual cortex.

## 110 4 Conclusion

111 We have shown that classic Gestalt phenomena—including closure, similarity grouping, good contin-  
 112 uation, and figure-ground segregation—can all be understood as expressions of a single mechanism:  
 113 feedback-driven prior integration. By reusing the same feedback pathways that support learning, our  
 114 theory-grounded algorithm *Prior-Guided Drift Diffusion* (PGDD) enables standard neural networks to  
 115 reproduce both the perceptual effects and the feedback dependence observed in biological vision. This  
 116 unification suggests that Gestalt principles are not independent heuristics but emergent consequences  
 117 of learned priors integrated at inference time. The approach highlights that perceptual organization,  
 118 long considered a gap between natural and artificial vision, can be realized without new architectures  
 119 or additional training. More broadly, our results point toward inference-time computations inspired  
 120 by human and animal brain, where feedback not only enables learning but also actively shapes  
 121 perception. Such dual use of feedback suggests a bridge between perceptual inference and continual  
 122 learning, where the same mechanism that guides moment-to-moment interpretation may also support  
 123 adaptation and knowledge accumulation across experience.

## References

- 124
- 125 [1] Anja Wannig, Liviu Stanisor, and Pieter R. Roelfsema. Automatic spread of attentional response  
126 modulation along gestalt criteria in primary visual cortex. *Nature Neuroscience*, 14(10):1243–  
127 1244, 2011.
- 128 [2] E. Leeuwenberg and P. A. van der Helm. Structural information theory: Towards a theory of  
129 perceptual organization. *Psychological Review*, 98(2):267–285, 1991.
- 130 [3] V. Froyen, J. Feldman, and M. Singh. Bayesian hierarchical grouping: Perceptual grouping as  
131 mixture estimation. *Psychological Review*, 122(4):575–597, 2015.
- 132 [4] Rajesh PN Rao and Dana H. Ballard. Predictive coding in the visual cortex: a functional  
133 interpretation of some extra-classical receptive-field effects. *Nature Neuroscience*, 2(1):79–87,  
134 1999.
- 135 [5] Zhaoyang Pang, Callum Biggs O’May, Bhavin Choksi, and Rufin VanRullen. Predictive coding  
136 feedback results in perceived illusory contours in a recurrent neural network. *Neural Networks*,  
137 144:164–175, 2021. doi: 10.1016/j.neunet.2021.08.024.
- 138 [6] Daniel L. K. Yamins, Ha Hong, Charles F. Cadieu, Ethan A. Solomon, Darren Seibert, and  
139 James J. DiCarlo. Performance-optimized hierarchical models predict neural responses in higher  
140 visual cortex. *Proceedings of the National Academy of Sciences*, 111(23):8619–8624, May  
141 2014. ISSN 1091-6490. doi: 10.1073/pnas.1403112111.
- 142 [7] Seyed-Mahdi Khaligh-Razavi and Nikolaus Kriegeskorte. Deep supervised, but not unsuper-  
143 vised, models may explain IT cortical representation. *PLoS Comput. Biol.*, 10(11):e1003915,  
144 November 2014.
- 145 [8] N. Baker, H. Lu, G. Erlikhman, and P. J. Kellman. Deep convolutional networks do not classify  
146 based on global object shape. *PLoS Computational Biology*, 14(12):e1006613, 2018. doi:  
147 10.1371/journal.pcbi.1006613.
- 148 [9] N. Baker, P. J. Kellman, G. Erlikhman, and H. Lu. Deep convolutional networks do not perceive  
149 illusory contours. In *Proceedings of the 40th Annual Conference of the Cognitive Science  
150 Society*, Madison, WI, 2018.
- 151 [10] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and  
152 Wieland Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias  
153 improves accuracy and robustness. In *International Conference on Learning Representations*,  
154 2019.
- 155 [11] Been Kim, Emily Reif, Martin Wattenberg, Samy Bengio, and Michael C. Mozer. Neural  
156 networks trained on natural scenes exhibit gestalt closure. *Computational Brain & Behavior*, 4  
157 (3):251–263, 2021.
- 158 [12] V. Biscione and J. S. Bowers. Mixed evidence for gestalt grouping in deep neural networks.  
159 *arXiv preprint arXiv:2203.07302*, 2022.
- 160 [13] V A Lamme. The neurophysiology of figure-ground segregation in primary visual cortex. *J.  
161 Neurosci.*, 15(2):1605–1615, February 1995.
- 162 [14] U. H. Schnabel, C. Bossens, J. A. M. Lorteije, et al. Figure-ground perception in the awake  
163 mouse and neuronal activity elicited by figure-ground stimuli in primary visual cortex. *Scientific  
164 Reports*, 8:17800, 2018. doi: 10.1038/s41598-018-36087-8. Received: 13 March 2018;  
165 Accepted: 09 November 2018; Published: 12 December 2018.
- 166 [15] Tai Sing Lee and My Nguyen. Dynamics of subjective contour formation in the early visual  
167 cortex. *Proceedings of the National Academy of Sciences*, 98(4):1907–1911, 2001. doi:  
168 10.1073/pnas.98.4.1907.
- 169 [16] D H Grosf, R M Shapley, and M J Hawken. Macaque V1 neurons can signal ‘illusory’ contours.  
170 *Nature*, 365(6446):550–552, October 1993.

- 171 [17] Tai Sing Lee and David Mumford. Hierarchical bayesian inference in the visual cortex. *JOSA*  
172 *A*, 20(7):1434–1448, 2003.
- 173 [18] Lauri Parkkonen, Jesper Andersson, Matti Hämäläinen, and Riitta Hari. Early visual brain areas  
174 reflect the percept of an ambiguous scene. *Proceedings of the National Academy of Sciences*,  
175 105(51):20500–20504, 2008. doi: 10.1073/pnas.0810966105.
- 176 [19] Alexandr Pak, Esther Ryu, Claudia Li, and Alexander A. Chubykin. Top-down feedback  
177 controls the cortical representation of illusory contours in mouse primary visual cortex. *The*  
178 *Journal of Neuroscience*, 40(3):648–660, December 2019. ISSN 1529-2401.
- 179 [20] H. Shin, M. B. Ogando, L. Abdeladim, et al. Recurrent pattern completion drives the neocortical  
180 representation of sensory inference. *bioRxiv*, 2023.
- 181 [21] Shibani Santurkar, Dimitris Tsipras, Brandon Tran, Andrew Ilyas, Logan Engstrom, and  
182 Aleksander Madry. Image synthesis with a single (robust) classifier. 2019.
- 183 [22] Harris Drucker and Yann Le Cun. Improving generalization performance using double back-  
184 propagation. In *Advances in Neural Information Processing Systems*, volume 5, pages 145–151,  
185 1992.
- 186 [23] Andrew Slavin Ross and Finale Doshi-Velez. Improving neural network generalization by  
187 reducing input gradient. In *Proceedings of the AAAI Conference on Artificial Intelligence*,  
188 volume 32, 2018.
- 189 [24] Judy Hoffman, Avinash Ravichandran, and Oncel Tuzel. Robust learning with jacobian reg-  
190 ularization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
191 *Recognition*, pages 4349–4357, 2019.
- 192 [25] Roman Novak, Yasaman Bahri, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-  
193 Dickstein. Sensitivity and generalization in neural networks: an empirical study. In *International*  
194 *Conference on Learning Representations*, 2018.
- 195 [26] Moustapha Cisse, Piotr Bojanowski, Edouard Grave, Yann Dauphin, and Nicolas Usunier.  
196 Parseval networks: Improving robustness to adversarial examples. In *Proceedings of the 34th*  
197 *International Conference on Machine Learning*, pages 854–863, 2017.
- 198 [27] Samet Oymak and Mahdi Soltanolkotabi. Towards moderate overparameterization: Global  
199 convergence guarantees for training shallow neural networks. In *Advances in Neural Information*  
200 *Processing Systems*, volume 32, 2019.
- 201 [28] Levent Sagun, Utku Evci, V Ugur Guney, Yann Dauphin, and Léon Bottou. Empirical analysis  
202 of the hessian of over-parametrized neural networks. *arXiv preprint arXiv:1706.04454*, 2018.
- 203 [29] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda B. Viégas, and Martin Wattenberg. Smooth-  
204 grad: removing noise by adding noise. In *Proceedings of the 34th International Conference on*  
205 *Machine Learning Workshops*, 2017.
- 206 [30] Jonas Kubilius, Martin Schrimpf, Ha Hong, Najib J. Majaj, Rishi Rajalingham,  
207 Elias B. Issa, Kohitij Kar, Pouya Bashivan, Jonathan Prescott-Roy, Kailyn Schmidt,  
208 Aran Nayebi, Daniel Bear, Daniel L. K. Yamins, and James J. DiCarlo. Brain-  
209 Like Object Recognition with High-Performing Shallow Recurrent ANNs. In  
210 H. Wallach, H. Larochelle, A. Beygelzimer, F. D’Alché-Buc, E. Fox, and R. Gar-  
211 nett, editors, *Neural Information Processing Systems (NeurIPS)*, pages 12785—  
212 12796. Curran Associates, Inc., 2019. URL [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/9441-brain-like-object-recognition-with-high-performing-shallow-recurrent-anns)  
213 [9441-brain-like-object-recognition-with-high-performing-shallow-recurrent-anns](http://papers.nips.cc/paper/9441-brain-like-object-recognition-with-high-performing-shallow-recurrent-anns).
- 214 [31] Kuan Han, Haiguang Wen, Yizhen Zhang, Di Fu, Eugenio Culurciello, and Zhongming Liu.  
215 Deep predictive coding network with local recurrent processing for object recognition, 2018.  
216 URL <https://arxiv.org/abs/1805.07526>.

217 **A Supplementary Material**

218 **A.1 Theoretical intuition: PGDD as Implicit Denoising**

219 The Prior-Guided Drift Diffusion (PGDD) algorithm operates through a mathematically principled  
 220 mechanism that leverages implicit denoising capabilities embedded in robust neural networks. Here  
 221 we provide the theoretical justification for why PGDD successfully accesses learned statistical priors.

222 **Gradient Structure in PGDD** Given the PGDD objective:

$$\mathcal{L}_{\text{PGDD}} = \|r(\mathbf{x}) - r(\mathbf{x} + \boldsymbol{\epsilon})\|^2 \quad (3)$$

223 where  $r(\cdot)$  represents network activations at layer  $\ell$  and  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ , the gradient with respect to  
 224 input  $\mathbf{x}$  becomes:

$$\nabla_{\mathbf{x}} \mathcal{L}_{\text{PGDD}} = 2\mathbf{J}(\mathbf{x})^T (\mathbf{r}(\mathbf{x}) - \mathbf{r}(\mathbf{x} + \boldsymbol{\epsilon})) \quad (4)$$

225 where  $\mathbf{J}(\mathbf{x}) = \nabla_{\mathbf{x}} r(\mathbf{x})$  is the Jacobian of representations with respect to input.

226 Using the linearization  $r(\mathbf{x} + \boldsymbol{\epsilon}) \approx r(\mathbf{x}) + \mathbf{J}(\mathbf{x})\boldsymbol{\epsilon}$ , this simplifies to:

$$\nabla_{\mathbf{x}} \mathcal{L}_{\text{PGDD}} \approx 2\mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x})\boldsymbol{\epsilon} = 2\mathbf{J}^T \mathbf{J} \boldsymbol{\epsilon} \quad (5)$$

227 **Low-Rank Structure and Denoising** Robust neural networks are known to exhibit reduced Jacobian  
 228 norms and improved local Lipschitz constants [22–24]. Adversarial training can be understood as  
 229 operator-norm regularization on  $\mathbf{J}$  [25, 26], which reshapes the singular value spectrum. Empirical  
 230 studies report that most singular values of  $\mathbf{J}$  are strongly suppressed while a few dominant modes  
 231 remain [27, 28]. This produces a low *effective* rank structure: sensitivity is concentrated in a small  
 232 discriminative subspace while class-insensitive directions are attenuated.

233 For a low-rank Jacobian  $\mathbf{J}$ , the operator  $\mathbf{J}^T \mathbf{J}$  inherits this structure and acts as a projection onto the  
 234 span of the network’s learned feature directions. When applied to noise  $\boldsymbol{\epsilon}$ :

$$\mathbf{J}^T \mathbf{J} \boldsymbol{\epsilon} = \mathbf{J}^T \mathbf{J} (\boldsymbol{\epsilon}_{\parallel} + \boldsymbol{\epsilon}_{\perp}) \quad (6)$$

235 where  $\boldsymbol{\epsilon}_{\parallel}$  lies in the span of learned features and  $\boldsymbol{\epsilon}_{\perp}$  is orthogonal to this span. Due to the low-rank  
 236 structure:

$$\mathbf{J}^T \mathbf{J} \boldsymbol{\epsilon} \approx \mathbf{J}^T \mathbf{J} \boldsymbol{\epsilon}_{\parallel} \quad (7)$$

237 This selective filtering preserves class-relevant perturbations while suppressing orthogonal noise  
 238 components, effectively implementing a learned denoising operator.

239 **Connection to Score-Based Models** The denoising operator  $\mathbf{J}^T \mathbf{J}$  in PGDD parallels the score  
 240 function  $\nabla_{\mathbf{x}} \log p(\mathbf{x})$  (estimated by denoiser autoencoders) used in diffusion models. Both operators  
 241 guide iterative refinement toward higher-probability regions of the learned data distribution. However,  
 242 unlike explicit generative models, PGDD accesses this denoising capability as an emergent property  
 243 of robust classification training.

244 The iterative update rule:

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \eta \mathbf{J}^T \mathbf{J} \boldsymbol{\epsilon} + \boldsymbol{\xi}_t \quad (8)$$

245 where  $\boldsymbol{\xi}_t \sim \mathcal{N}(0, \sigma_{\text{diff}}^2 \mathbf{I})$  provides stochastic exploration, implements a form of Langevin dynamics  
 246 guided by the implicit score function encoded in the robust classifier’s gradient structure. This  
 247 theoretical framework explains why PGDD can successfully extract learned priors from networks  
 248 trained solely for classification, without requiring additional generative objectives or architectural  
 249 modifications.

250 **A.2 PGDD algorithm**

---

**Algorithm 1** Prior-Guided Drift Diffusion Objective

---

```

1: Input: Image  $x_{\text{input}}$ , model  $f$ , target layer  $\ell$ , constraint  $\epsilon$ , step size  $\alpha$ , noise ratio  $\sigma$ , iterations  $T$ 
2: Output: Refined representations  $\{x_t\}_{t=0}^T$ 
3: // Step 1: Feedforward pass
4:  $x_0 \leftarrow \text{normalize}(x_{\text{input}})$ 
5:  $f_\ell \leftarrow \text{extract\_layers}(f, \ell)$  {Extract model up to layer  $\ell$ }
6:  $x_{\text{noisy}} \leftarrow x_0 + \sigma \cdot \mathcal{N}(0, I)$ 
7:  $r_{\text{anti-target}} \leftarrow f_\ell(x_{\text{noisy}})$  {Generate noisy reference representation}
8: for  $t = 1$  to  $T$  do
9:   // Step 2: Inference objective selection
10:   $\text{anti-target} \leftarrow r_{\text{anti-target}}$  {Use noisy reference as target}
11:  // Step 3: Feedback error propagation
12:   $h_t \leftarrow f_\ell(x_{t-1})$  {Forward pass through target layers}
13:   $\mathcal{L}_t \leftarrow \|h_t - r_{\text{anti-target}}\|^2$  {MSE loss in representation space}
14:   $g_t \leftarrow \nabla_{x_{t-1}} \mathcal{L}_t$  {Gradient via feedback pathways}
15:  // Step 4: Constrained activation update
16:   $\tilde{g}_t \leftarrow \alpha \cdot g_t / (\|g_t\| + 1e-10)$  {Normalize gradient}
17:   $\eta_t \leftarrow \text{diffusion\_noise\_ratio} \cdot \mathcal{N}(0, I)$  {Add stochastic noise}
18:   $x'_t \leftarrow x_{t-1} + \tilde{g}_t + \eta_t$  {Move away from representation of noisy input (anti-target)}
19:   $x_t \leftarrow \text{project}(x'_t, x_0, \epsilon)$  {Enforce  $\|x_t - x_0\|_\infty \leq \epsilon$ }
20:  // Step 5: Iteration control
21:  {Continue to next iteration}
22: end for
23:
24: return  $\{x_t\}_{t=0}^T$ 

```

---

251 **A.3 PGDD parameters**

252 **Note:** These parameters represent working values for the demonstrations shown. PGDD operates  
253 effectively over broader parameter ranges not detailed here.

Table 1: PGDD parameters for different Gestalt principle instances

Phenomenon	Layer	Noise $\sigma^2$	Iterations	Step Size (update rate)
Closure: Kanizsa figures	Last	0.1	50	0.5
Similarity grouping	Layer 3	0.1	50	1.0
Good continuation	Layer 3	0.1	50	0.4
Figure-ground	Layer 3	0.5	100	0.8
Rubin’s vase	Last	0.5	100	0.8

254 **A.4 sPGDD: probabilistic variant of PGDD for standard neural networks**

255 While PGDD works optimally with networks trained for robustness to adversarial perturbations [21],  
 256 PGDD can still be extended to standard networks trained without explicit robustness objectives,  
 257 we developed a probabilistic variant called smooth Prior-Guided Drift Diffusion (sPGDD). sPGDD  
 258 leverages averaging gradients over multiple noise-corrupted versions of an input [29]. At each  
 259 iteration, sPGDD generates multiple noise-corrupted versions of the current state, computes gradients  
 for each sample, and averages these gradients to update the input activation.

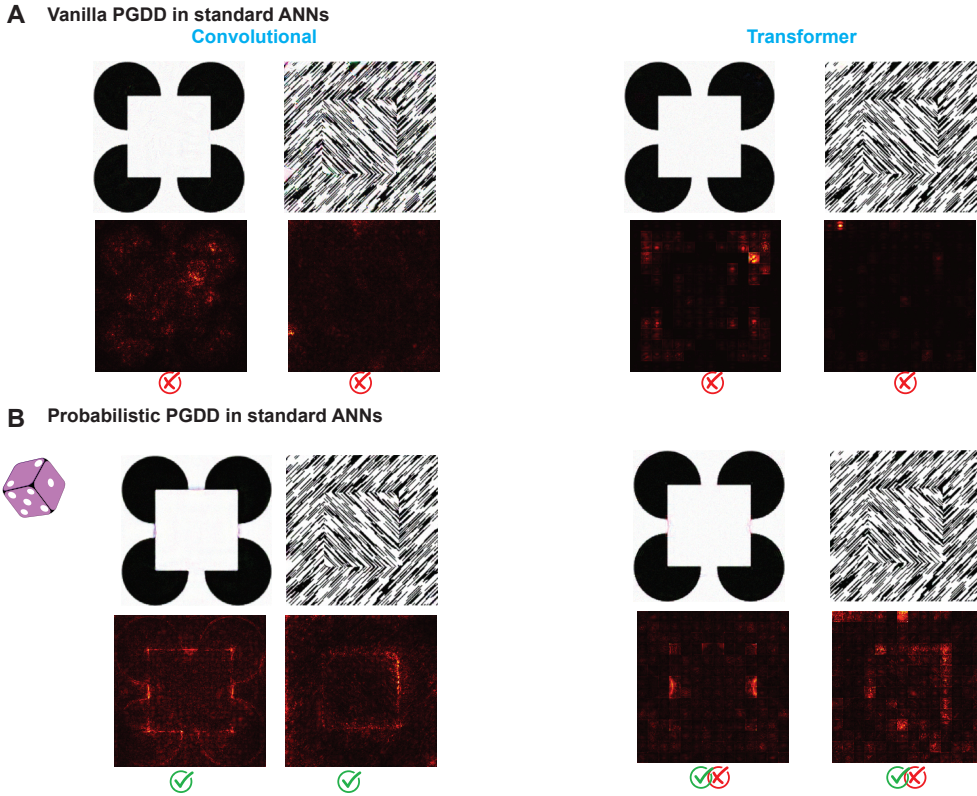


Figure S1: sPGDD extends PGDD to standard neural networks through probabilistic inference. (A) Standard generative inference fails in conventional networks: Both convolutional (ResNet-50) and transformer (ViT-Large) architectures trained without robustness show minimal illusory contour formation when using direct PGDD on Kanizsa squares and figure-ground stimuli. (B) Probabilistic PGDD (or smoothed PGDD) succeeds through gradient averaging [29]: The same networks generate clear illusory contours and figure-ground effects when using sPGDD with gradient averaging over multiple noise samples.

260

261 **A.5 End-to-end training of feedback as in RNNs and predictive coding**

262 To demonstrate the specificity of PGDD, we investigate alternative frameworks that incorporate  
 263 feedback processing but use separate feedback for learning and inference. We tested two prominent  
 264 architectures that explicitly include recurrent or feedback connections trained end-to-end for visual  
 265 recognition: CORnet-S [30] and PredNetBpE [31].

266 CORnet-S implements recurrent processing within cortical areas through fixed recurrent weights  
 267 during inference, effectively extending processing depth but not directly leveraging the same pathways  
 268 used for error propagation during learning. PredNetBpE implements predictive coding principles  
 269 with feedback connections that carry top-down predictions and bottom-up prediction errors during  
 270 both training and inference. Both models were trained on ImageNet for object recognition and tested  
 271 on the same Gestalt stimuli used to evaluate PGDD. Despite incorporating feedback mechanisms,  
 272 neither model reliably reproduced the range of Gestalt phenomena observed with PGDD, highlighting  
 273 that end-to-end training of recurrent feedback (through backpropagation through time) is insufficient  
 to account for these perceptual and neural phenomenon.

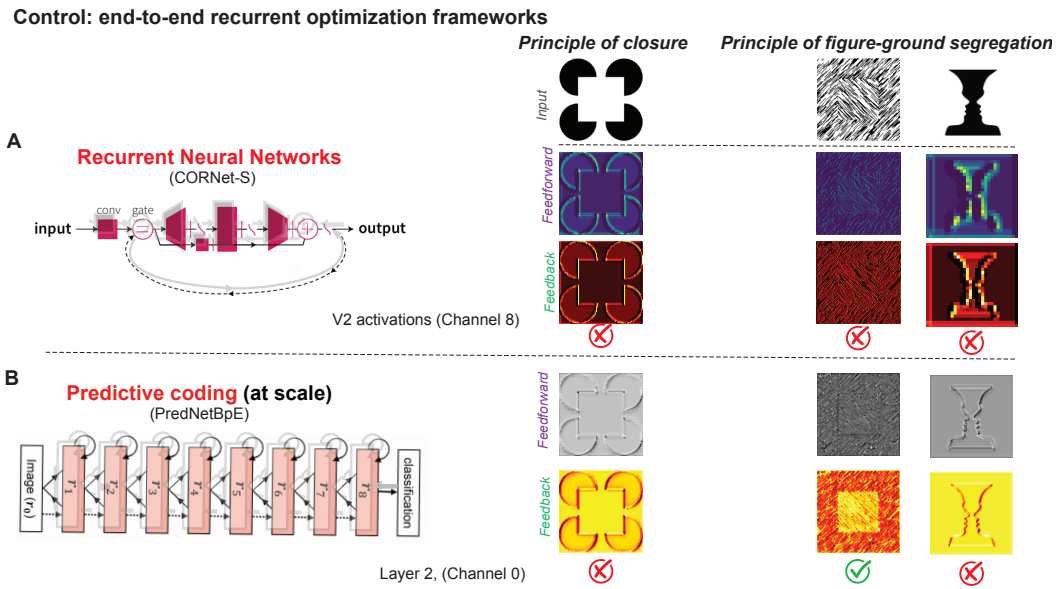


Figure S2: Alternative recurrent frameworks show limited Gestalt processing compared to PGDD. (A) CORnet-S, a biologically-inspired recurrent neural network, fails to generate clear closure effects for Kanizsa figures (red X) and shows weak figure-ground segregation despite its recurrent architecture. (B) PredNetBpE, implementing predictive coding at scale, captures some figure-ground segregation (green checkmark) but fails at closure generation.

274