
Common visual learning constraints in transformers and newborn brains: Evidence from line drawings

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A core goal in artificial intelligence (AI) is to build machines that learn like brains.
2 Many AI systems, including convolutional neural networks (CNNs) and vision
3 transformers (ViTs), rival human adults on visual recognition tasks. But, do these
4 AI systems actually learn like brains? If so, AI systems should produce the same
5 learning outcomes as brains when trained with the same data. Here, we tested
6 whether AI systems learn the same object recognition skills as newborn chicks
7 when trained in the same visual environments as chicks. We performed digital twin
8 studies of prior controlled-rearing experiments, evaluating whether CNNs and ViTs
9 produce the same pattern of successes and failures as chicks. When ViTs were
10 equipped with a biologically inspired temporal learning objective, the ViTs showed
11 the same learning patterns as chicks: both learned object recognition when reared
12 with normal objects, but failed to learn object recognition when reared with line
13 drawings. Conversely, when CNNs were equipped with the same temporal learning
14 objective, the CNNs showed a different pattern from chicks: CNNs learned object
15 recognition whether exposed to normal objects or line drawings. These results show
16 that transformers can be accurate image-computable models of visual learning.

17 1 Introduction

18 A major scientific and engineering goal is to build machines that learn like brains. For science, this
19 would provide working models for simulating how brains learn to perceive and understand the world.
20 For engineering, this would provide systems that learn with the same power and efficiency as brains.
21 The past decade has produced dozens of success stories in which machines match or exceed the
22 abilities of humans. For example, convolutional neural networks (CNNs) and vision transformers
23 (ViTs) can achieve high levels of performance on a range of tasks, including object recognition [1, 2],
24 action recognition [3], scene perception [4], object segmentation [5], optic flow perception [6], and
25 navigation [7].

26 But, do AI systems actually learn like brains (i.e., produce the same learning outcomes when trained
27 with the same data)? As many researchers have pointed out [8–11], the training regimes faced by
28 animals and machines differ radically. AI systems are typically trained on massive datasets (e.g.,
29 millions of images and videos across thousands of object categories and environments), whereas
30 animals spend their postnatal lives in one environment surrounded by a handful of objects and
31 caregivers. On face, there seems to be a massive mismatch between the volume and variety of
32 training data needed by machines versus animals. Accordingly, AI systems are often regarded as
33 data hungry systems, “gorging on hundreds of terabytes of data,” whereas brains are thought to be
34 “efficient and even elegant systems that operate with small amounts of information” [12]. From this
35 perspective, learning in brains and machines is nothing alike.

36 The view that AI systems are data hungry rests on the assumption that the visual experiences of new-
37 borns are impoverished and noisy, a “blooming, buzzing confusion” [13]. However, recent work from

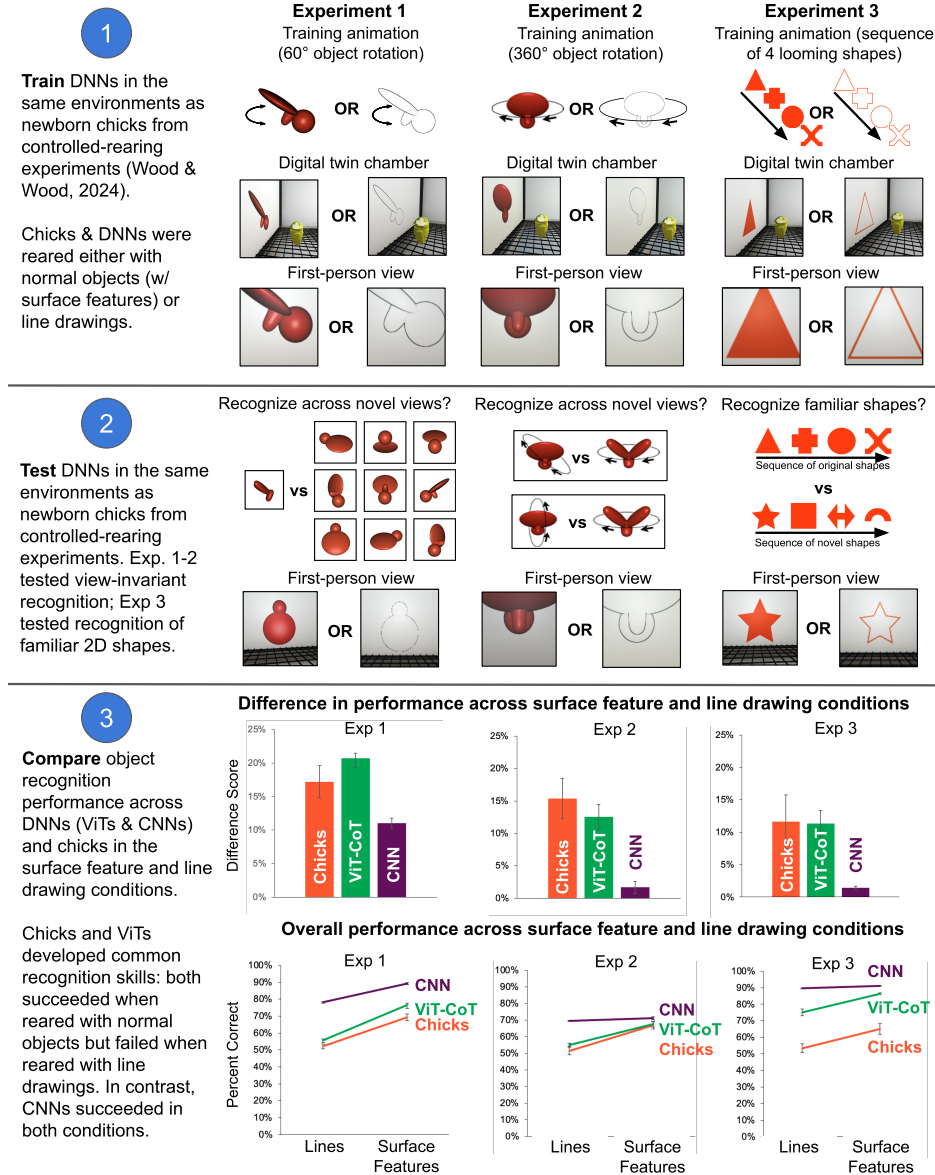


Figure 1: Design and Results. (1) Deep neural networks (DNNs) and newborn chicks were reared in the same visual environments, containing either normal objects or line drawings. (2) DNNs and chicks were tested with the same object recognition tasks. (3) Chicks and vision transformers (ViTs) showed common patterns of development: both learned object recognition when reared with normal objects, and both failed to learn object recognition when reared with line drawings. In contrast, convolutional neural networks (CNNs) learned object recognition in both conditions.

38 developmental psychology suggests that the opposite is true: The first-person views acquired by babies
 39 are temporally and spatially rich [14]. By moving their bodies and heads, babies produce large num-
 40 bers of diverse, high-quality object views that are well suited for learning[10, 15]. In fact, when first-
 41 person views from babies are used to train CNNs and ViTs, the models learn core visual skills [9, 11].
 42 Thus, the gap between human and machine vision might not be as great as previously thought [16].

43 Human infants acquire a rich visual diet filled with many objects, people, and environments.
 44 However, biological visual systems can learn effectively even in impoverished worlds. For example,

45 newborn chicks learn object perception in worlds that contain just a single object [17–20]. Can
46 CNNs and ViTs learn in the same impoverished environments faced by newborn animals?

47 Digital twin studies [21] are designed to tackle this question, by raising animals and machines in the
48 same environments and testing them with the same tasks. Researchers control and match the training
49 data from which brains and machines learn, allowing for direct comparison of learning outcomes.
50 Prior digital twin studies show that AI algorithms show common learning successes as newborn
51 animals. When CNNs [22, 23] or ViTs [8] are trained on first-person views of agents exploring
52 virtual animal chambers that mimic the rearing conditions of chicks, CNNs and ViTs learn the same
53 object recognition skills as chicks. These results contradict the view that AI algorithms are more
54 data hungry than brains.

55 The discovery that both CNNs and ViTs can learn effectively in the impoverished environments faced
56 by newborn animals raises a new challenge: how do we distinguish between these model classes? On
57 one hand, CNNs might be the more accurate model class because CNNs and newborn visual systems
58 are both hierarchically and retinotopically organized [24, 25]. On the other hand, the visual system’s
59 receptive field structure could be an emergent property of even more foundational (and generic)
60 learning machinery. For instance, fully connected neural networks learn convolutional structures
61 when trained on data with non-Gaussian, higher-order local structure [26]. During prenatal develop-
62 ment, brains are shaped by spontaneous retinal waves, which have a non-Gaussian, higher-order local
63 structure [27]. Thus, brains could learn a hierarchical and retinotopic organization during prenatal
64 development, powered by more generic learning mechanisms. If so, the core learning machinery
65 driving visual intelligence would not be CNN-like; rather, it might be more like a transformer. Indeed,
66 ViTs learn CNN-like receptive field structures when trained on natural images [28]. The core com-
67 putations in transformers also closely match those in the neuron–astrocyte network in the brain [29],
68 raising the possibility that ViTs are the more accurate model of the core learning algorithm in brains.

69 We tested this hypothesis by evaluating whether CNNs and ViTs show the same successes *and*
70 *failures* as newborn chicks. Newborn chicks learn better from some experiences than others, so by
71 examining whether CNNs and ViTs show the same pattern of successes and failures as newborn
72 chicks across studies, we can measure which model class learns more like brains. We focused on
73 visual learning from normal objects versus line drawings (objects lacking surface features, Fig. 1).
74 If a chick’s visual environment contains normal objects with surface features, then chicks learn to
75 recognize objects across familiar and novel viewpoints [17]. But, if a chick’s environment contains
76 line drawings, then chicks fail to develop object recognition [30]. For newborn brains, a visual diet
77 of line drawings is insufficient to learn object recognition.

78 **Line drawings for studying vision.** Line drawings have been used for decades to study object
79 recognition. Many studies show that human adults can readily recognize objects depicted in line
80 drawings (e.g., [31–34]). This ability develops rapidly. Infants show enhanced attention to lines
81 that depict corners and edges in the first year of life [35], and young children use lines to depict
82 objects in their earliest attempts to draw the world [36]. Humans have used line drawings to depict
83 scenes since prehistoric times [37, 38]. There is also evidence that nonhuman animals can recognize
84 line drawings, including chimpanzees [39, 40] and pigeons [41].

85 None of these studies, however, tested humans or animals at the beginning of life. All of the subjects
86 had already acquired months to years of visual experience with real-world objects before they were
87 tested. To explore whether newborn brains can recognize line drawings, Wood and Wood [30] used
88 controlled rearing. The researchers raised newborn chicks in automated controlled-rearing chambers
89 that contained a single object, then tested the chicks’ ability to recognize that object across novel
90 viewpoints. When chicks were reared with an object that had surface features, the chicks developed
91 view-invariant object recognition. However, when chicks were reared with a line drawing of an
92 object, the chicks failed to develop object recognition. Do CNNs and ViTs show this same learning
93 pattern? We address this question through digital twin experiments, raising CNNs and ViTs in the
94 same visual environments as chicks and testing them with the same tasks.

95 2 Methods

96 **Architecture.** We used two architectures (CNNs and ViTs) because both are high-performing model
97 classes on a range of visual recognition tasks [42–44] and because the models differ in terms of their
98 hardcoded inductive biases. CNNs have a strong spatial bias. The convolutional operation reflects

99 the spatial structure of natural images, allowing CNNs to generalize well from small datasets and
100 learn useful feature hierarchies that capture the structure of visual images [45, 46]. Conversely, ViTs
101 are generic learning algorithms that do not have hardcoded knowledge about objects or space [43].
102 Instead, ViTs learn through flexible (learned) allocation of attention that does not assume any spatial
103 (or object) structure.

104 **Objective Function.** For each experiment, we performed comparisons between CNNs and ViTs that
105 had the same temporal learning objective function. Based on decades of empirical and theoretical
106 work in neuroscience, we hypothesized that unsupervised temporal learning (UTL) drives visual
107 development in the brain [47–52]. According to UTL models, brains build object representations
108 by adapting to the spatiotemporal statistics of the animal’s visual environment. The key assumption
109 underlying UTL is that distal scene variables (e.g., curvature, depth, orientation, texture, shape) vary
110 slowly over time in natural visual environments. Thus, in principle, brains could learn distal scene
111 variables by encoding statistical regularities across successive changes in proximal retinal images
112 (see Appendix A.4 for details).

113 Our models were initially untrained (no pre-training), and during training, the models were trained
114 on simulated first-person visual experiences from chicks. Like the chicks, the models’ visual diet
115 was limited to a single object in a controlled-rearing chamber.

116 **Training Data** We used behavioral benchmarks from newborn chicks because chicks can be raised
117 in strictly controlled environments from the onset of vision, providing strict control of all visual
118 experiences (training data) acquired by the animal [53]. This control over training data is essential for
119 directly comparing learning across animals and machines. Chicks can also inform our understanding
120 of human vision because avian brains have similar cells and circuitry as mammalian brains [54–56],
121 as well as a similar large-scale organization, including a hierarchy of sensory information processing,
122 hippocampus regions, and associative areas.

123 To simulate the visual experiences of newborn chicks, we created realistic digital twins of the
124 controlled-rearing chambers in a video game engine (Unity 3D). Then, we simulated the visual
125 diet available in the chick’s environment by recording the first-person images acquired by an agent
126 moving through the virtual chambers. We collected 80,000 first-person images in each of the rearing
127 conditions and used those images to train the CNNs and ViTs (see Appendix section A.2 for details).

128 3 Results

129 3.1 Experiment 1: 60° object rotation experience

130 In Experiment 1 (Fig. 1, *left*), we focused on the view-invariant object recognition task and
131 data reported in Wood [17] and Wood & Wood [30]. Newborn chicks were hatched in darkness,
132 then raised singly in automated controlled-rearing chambers that measured each chick’s behavior
133 continuously (24/7) during the first two weeks of life. The chambers were equipped with two display
134 walls (LCD monitors) for displaying object stimuli. The chambers did not contain any objects other
135 than the virtual objects projected on the display walls, providing control over all object experiences
136 acquired by the animal from the onset of vision.

137 During the training phase, chicks were reared in an environment containing a single virtual object rotat-
138 ing through a 60° viewpoint range. This virtual object was the only object in the chick’s environment.
139 The chicks were raised in this environment for 1 week, allowing the critical period on filial imprinting
140 to close. The chicks were raised and tested with either line drawings or objects with surface features.

141 During the test phase, the chicks were tested on their ability to recognize the imprinted object across
142 12 in-depth viewpoint changes. On each test trial, the imprinted object appeared on one display wall
143 and an unfamiliar object appeared on the opposite display wall. Test trials were scored as correct when
144 the chicks spent a greater proportion of time with their imprinted object and incorrect when the chicks
145 spent a greater proportion of time with the unfamiliar object. The viewpoint changes introduced large,
146 novel, and complex changes in the object’s appearance. Nevertheless, as shown in Fig. 1 (*Panel 3, left*),
147 the chicks reared with surface-feature objects successfully recognized their imprinted object across
148 the novel viewpoints. From a visual diet of a single object, chicks can learn view-invariant object
149 representations. In contrast, when chicks were reared with line drawings of that same object, the chicks
150 never learned to recognize objects. The chicks reared with the line drawings performed at chance level,
151 despite acquiring over 100 hours of visual experience with the line drawings during the training phase.

152 To compare learning across chicks, CNNs, and ViTs, we performed matching controlled-rearing
153 experiments on CNNs and ViTs (Fig. 1, *Panels 1 & 2*). We created digital twins of the controlled-
154 rearing chambers, then simulated the visual diet in those chambers and used those simulated data
155 streams to train CNNs and ViTs. We then tested the models with the same stimuli used to test the
156 chicks. The chicks and models were trained in the same visual environment and tested on the same
157 task, allowing for direct comparison of their learning outcomes.

158 Fig. 1 (*Panel 3, left*), shows the performance of CNNs (SimCLR-CLTT) and ViTs (ViT-CoT) in
159 the surface feature and line drawing conditions. SimCLR-CLTT succeeded at the task, learning
160 view-invariant object representations in both conditions. In contrast, ViT-CoT showed the same
161 learning pattern as chicks: ViT-CoT succeeded when learning from normal objects, but failed when
162 learning from line drawings.

163 **3.2 Experiment 2: 360° of object rotation experience**

164 To validate this conclusion under different conditions, we performed a second digital twin experiment
165 of prior controlled-rearing studies [30, 57](Fig. 1, *middle*). Rather than presenting the objects from
166 a 60° viewpoint range, the objects moved through a 360° viewpoint range, completing an in-depth
167 rotation every 15 seconds. The chicks, CNNs, and ViTs were thus exposed to six times as many
168 unique views of the object during the training phase. In the test phase, we tested whether the models
169 could recognize the imprinted object across novel viewpoints, by rotating the object around novel
170 axes of rotation (Fig. 1, *Panels 1 & 2*).

171 As shown in Fig. 1 (*Panel 3, middle*), when chicks were reared with an object with surface features,
172 the chicks built view-invariant representations that generalized across large, novel, and complex
173 changes in the object’s appearance [57]. When chicks were reared with line drawings, they performed
174 at chance level, despite acquiring over 100 hours of visual experience with the line drawings during
175 the training phase [30]. Fig. 1 (*Panel 3, middle*) shows the performance of SimCLR-CLTT and
176 ViT-CoT in the surface feature and line drawing conditions. Again, SimCLR-CLTT succeeded on the
177 task, learning view-invariant object representations in both conditions. In contrast, ViT-CoT showed
178 the same learning pattern as the chicks. ViT-CoT succeeded when learning from normal objects,
179 but failed when learning from line drawings.

180 **3.3 Experiment 3: Looming 2D shapes**

181 To validate our results with different object stimuli, we performed a third digital twin experiment of
182 prior controlled-rearing studies [30, 58]. These studies used simple two-dimensional objects, rather
183 than complex three-dimensional objects. During the training phase, the chicks were presented with
184 a sequence of four looming shapes (Fig. 1, *right*). During the test phase, the chicks were tested on
185 their ability to distinguish familiar shapes from novel shapes.

186 When chicks were reared with a sequence of shapes containing surface features, they reliably
187 distinguished familiar from novel shapes [58]. In contrast, when reared with a sequence of line
188 drawing shapes, the chicks failed to distinguish familiar from novel shapes [30]. Fig. 1 (*Panel 3,*
189 *right*) shows the performance of SimCLR-CLTT and ViT-CoT in the surface feature and line drawing
190 conditions. SimCLR-CLTT performed equally well in the surface features and lines conditions.
191 Conversely, ViT-CoT, like the chicks, showed impaired recognition when learning from line drawings.

192 **4 Conclusion**

193 Do AI systems learn like brains? We trained CNNs and ViTs on simulated visual experiences from
194 newborn chicks, and found that temporal learning ViTs (ViT-CoT) showed the same learning patterns
195 as chicks. Both ViT-CoT and chicks learned object recognition when reared with normal objects,
196 but failed to learn object recognition when reared with line drawings. Conversely, CNNs equipped
197 with the same temporal learning objective as the ViTs (SimCLR-CLTT) did not show this pattern:
198 SimCLR-CLTT learned object recognition from both normal objects and line drawings. Appendix
199 A.1 contains additional experiments using alternative architectures and objective functions, and
200 Appendix A.7 contains a more detailed discussion of the limitations and broader impacts of this
201 work. Transformers, but not CNNs, showed the same visual learning pattern as chicks. We conclude
202 that transformers can be accurate image-computable models of visual learning in newborn brains.

203 **References**

- 204 [1] Priyank Jaini, Kevin Clark, and Robert Geirhos. Intriguing properties of generative classifiers.
205 *arXiv preprint arXiv:2309.16779*, 2023.
- 206 [2] Robert Geirhos, Kantharaju Narayanappa, Benjamin Mitzkus, Tizian Thieringer, Matthias
207 Bethge, Felix A Wichmann, and Wieland Brendel. Partial success in closing the gap between
208 human and machine vision. *Advances in Neural Information Processing Systems*, 34:23885–
209 23899, 2021.
- 210 [3] Anwaar Ulhaq, Naveed Akhtar, Ganna Pogrebna, and Ajmal Mian. Vision transformers for
211 action recognition: A survey. *arXiv preprint arXiv:2209.05700*, 2022.
- 212 [4] Aditya Jonnalagadda and Miguel Eckstein. A foveated vision-transformer model for scene
213 classification. *Journal of Vision*, 22(14):4440–4440, 2022.
- 214 [5] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson,
215 Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In
216 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026,
217 2023.
- 218 [6] Daniel M Bear, Kevin Feigelis, Honglin Chen, Wanhee Lee, Rahul Venkatesh, Klemen Kotar,
219 Alex Durango, and Daniel LK Yamins. Unifying (machine) vision via counterfactual world
220 modeling. *arXiv preprint arXiv:2306.01828*, 2023.
- 221 [7] Dhruv Shah, Ajay Sridhar, Nitish Dashora, Kyle Stachowicz, Kevin Black, Noriaki Hi-
222 rose, and Sergey Levine. Vint: A foundation model for visual navigation. *arXiv preprint*
223 *arXiv:2306.14846*, 2023.
- 224 [8] Lalit Pandey, Samantha Marie Waters Wood, and Justin Newell Wood. Are vision transformers
225 more data hungry than newborn visual systems? In *Thirty-seventh Conference on Neural*
226 *Information Processing Systems*, 2023.
- 227 [9] Chengxu Zhuang, Siming Yan, Aran Nayebi, Martin Schrimpf, Michael C. Frank, James J.
228 DiCarlo, and Daniel L. K. Yamins. Unsupervised neural network models of the ventral visual
229 stream. *Proceedings of the National Academy of Sciences*, 118(3):e2014196118, 2021. doi:
230 10.1073/pnas.2014196118.
- 231 [10] Saber Sheybani, Himanshu Hansaria, Justin Wood, Linda Smith, and Zoran Tiganj. Curriculum
232 learning with infant egocentric videos. In *Advances in Neural Information Processing Systems*,
233 volume 36, 2023.
- 234 [11] A Emin Orhan and Brenden M Lake. Learning high-level visual representations from a child’s
235 perspective without strong inductive biases. *Nature Machine Intelligence*, pages 1–13, 2024.
- 236 [12] Noam Chomsky, Ian Roberts, and Jeffrey Watumull. Noam chomsky: The false promise of
237 chatgpt. *The New York Times*, 8, 2023.
- 238 [13] W. James. *The Principles of Psychology*. Number v. 1 in American science series–Advanced
239 course. H. Holt, 1890.
- 240 [14] Linda B Smith, Swapnaa Jayaraman, Elizabeth Clerkin, and Chen Yu. The developing infant
241 creates a curriculum for statistical learning. *Trends in Cognitive Sciences*, 22(4):325–336, 2018.
- 242 [15] Sven Bambach, David Crandall, Linda Smith, and Chen Yu. Toddler-inspired visual object
243 learning. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- 244 [16] Justin N Wood. Artificial intelligence tackles the nature–nurture debate. *Nature Machine*
245 *Intelligence*, pages 1–2, 2024.
- 246 [17] Justin N Wood. Newborn chickens generate invariant object representations at the onset of visual
247 object experience. *Proceedings of the National Academy of Sciences*, 110(34):14000–14005,
248 2013.

- 249 [18] Justin N. Wood. Newly hatched chicks solve the visual binding problem. *Psychological Science*,
250 25(7):1475–1481, 2014.
- 251 [19] Justin N Wood and Samantha MW Wood. One-shot learning of view-invariant object represen-
252 tations in newborn chicks. *Cognition*, 199:104192, 2020.
- 253 [20] Samantha M. W. Wood and Justin N. Wood. One-shot object parsing in newborn chicks. *Journal*
254 *of Experimental Psychology. General*, 2021.
- 255 [21] Justin N Wood, Lalit Pandey, and Samantha MW Wood. Digital twin studies for reverse
256 engineering the origins of visual intelligence. *Annual Review of Vision Science*, 2024. In press.
- 257 [22] Lalit Pandey, Donsuk Lee, Samantha MW Wood, and Justin N Wood. Parallel development of
258 object recognition in newborn chicks and deep neural networks. Under review.
- 259 [23] Donsuk Lee, Pranav Gujarathi, and Justin N Wood. Controlled-rearing studies of newborn
260 chicks and deep neural networks. *arXiv preprint arXiv:2112.06106*, 2021.
- 261 [24] Michael J Arcaro and Margaret S Livingstone. On the relationship between maps and domains
262 in inferotemporal cortex. *Nature Reviews Neuroscience*, 22(9):573–583, 2021.
- 263 [25] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444,
264 2015.
- 265 [26] Alessandro Ingrosso and Sebastian Goldt. Data-driven emergence of convolutional structure
266 in neural networks. *Proceedings of the National Academy of Sciences*, 119(40):e2201854119,
267 2022.
- 268 [27] Mark V Albert, Adam Schnabel, and David J Field. Innate visual learning through spontaneous
269 activity patterns. *PLoS Computational Biology*, 4(8):e1000137, 2008.
- 270 [28] Muhammad Muzammal Naseer, Kanchana Ranasinghe, Salman H Khan, Munawar Hayat,
271 Fahad Shahbaz Khan, and Ming-Hsuan Yang. Intriguing properties of vision transformers.
272 *Advances in Neural Information Processing Systems*, 34:23296–23308, 2021.
- 273 [29] Leo Kozachkov, Ksenia V. Kastanenko, and Dmitry Krotov. Building transformers from neurons
274 and astrocytes. *Proceedings of the National Academy of Sciences*, 120(34):e2219150120, 2023.
275 doi: 10.1073/pnas.2219150120.
- 276 [30] Justin N Wood and Samantha MW Wood. The development of object recognition requires
277 experience with the surface features of objects. *Animals*, 14(2):284, 2024.
- 278 [31] Irving Biederman. Recognition-by-components: a theory of human image understanding.
279 *Psychological Review*, 94(2):115, 1987.
- 280 [32] Irving Biederman and Ginny Ju. Surface versus edge-based determinants of visual recognition.
281 *Cognitive Psychology*, 20(1):38–64, 1988.
- 282 [33] Alomit Ishai, Leslie Ungerleider, Alex Martin, and James Haxby. The representation of objects
283 in the human occipital and temporal cortex. *Journal of Cognitive Neuroscience*, 12 Suppl 2:
284 35–51, 11 2000. doi: 10.1162/089892900564055.
- 285 [34] Dirk B. Walther, Barry Chai, Eamon Caddigan, Diane M. Beck, and Li Fei-Fei. Simple line
286 drawings suffice for functional mri decoding of natural scene categories. *Proceedings of the*
287 *National Academy of Sciences*, 108(23):9661–9666, 2011. doi: 10.1073/pnas.1015666108.
- 288 [35] Albert Yonas and Martha E Arterberry. Infants perceive spatial structure specified by line
289 junctions. *Perception*, 23(12):1427–1435, 1994. doi: 10.1068/p231427. PMID: 7792132.
- 290 [36] Jacqueline Goodnow. *Children drawing*. Harvard University Press, 1977.
- 291 [37] Jean Clottes. Chauvet cave (ca. 30,000 bc), 2000.
- 292 [38] John M Kennedy and Abraham S Ross. Outline picture perception by the songe of papua.
293 *Perception*, 4(4):391–406, 1975.

- 294 [39] Shoji Itakura. Recognition of line-drawing representations by a chimpanzee (pan troglodytes).
295 *The Journal of General Psychology*, 121(3):189–197, 1994.
- 296 [40] Masayuki Tanaka. Recognition of pictorial representations by chimpanzees (pan troglodytes).
297 *Animal cognition*, 10:169–179, 2007.
- 298 [41] Edward A Wasserman, Joseph L Gagliardi, Brigitte R Cook, Kim Kirkpatrick-Steger, Suzette L
299 Astley, and Irving Biederman. The pigeon’s recognition of drawings of depth-rotated stimuli.
300 *Journal of Experimental Psychology: Animal Behavior Processes*, 22(2):205, 1996.
- 301 [42] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are
302 data-efficient learners for self-supervised video pre-training. *Advances in Neural Information*
303 *Processing Systems*, 35:10078–10093, 2022.
- 304 [43] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
305 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.
306 An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*
307 *arXiv:2010.11929*, 2020.
- 308 [44] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for im-
309 age recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
310 *Recognition*, pages 770–778, 2016.
- 311 [45] Yun-Hao Cao and Jianxin Wu. A random cnn sees objects: One inductive bias of cnn and its
312 applications. In *Proceedings Of The AAAI Conference On Artificial Intelligence*, volume 36,
313 pages 194–202, 2022.
- 314 [46] Shuying Liu and Weihong Deng. Very deep convolutional neural network based image clas-
315 sification using small training sample size. In *2015 3rd IAPR Asian conference on pattern*
316 *recognition (ACPR)*, pages 730–734. IEEE, 2015.
- 317 [47] James J DiCarlo, Davide Zoccolan, and Nicole C Rust. How does the brain solve visual object
318 recognition? *Neuron*, 73(3):415–434, 2012.
- 319 [48] Jacob Feldman and Patrice D Tremoulet. Individuation of visual objects over time. *Cognition*,
320 99(2):131–165, 2006.
- 321 [49] Peter Földiák. Learning invariance from transformation sequences. *Neural Computation*, 3(2):
322 194–200, 1991.
- 323 [50] Edmund T Rolls. Invariant visual object and face recognition: neural and computational bases,
324 and a model, visnet. *Frontiers in Computational Neuroscience*, 6:35, 2012.
- 325 [51] James V Stone. Learning perceptually salient visual parameters using spatiotemporal smooth-
326 ness constraints. *Neural Computation*, 8(7):1463–1492, 1996.
- 327 [52] Laurenz Wiskott and Terrence J Sejnowski. Slow feature analysis: Unsupervised learning of
328 invariances. *Neural Computation*, 14(4):715–770, 2002.
- 329 [53] Samantha MW Wood and Justin N Wood. A chicken model for studying the emergence of
330 invariant object recognition. *Frontiers in Neural Circuits*, 9:7, 2015.
- 331 [54] Onur Güntürkün and Thomas Bugnyar. Cognition without cortex. *Trends in Cognitive Sciences*,
332 20(4):291–303, 2016.
- 333 [55] Erich D Jarvis, Onur Güntürkün, Laura Bruce, András Csillag, Harvey Karten, Wayne Kuenzel,
334 Loreta Medina, George Paxinos, David J Perkel, Toru Shimizu, et al. Avian brains and a new
335 understanding of vertebrate brain evolution. *Nature Reviews Neuroscience*, 6(2):151–159, 2005.
- 336 [56] Harvey J Karten. Neocortical evolution: neuronal circuits arise independently of lamination.
337 *Current Biology*, 23(1):R12–R15, 2013.
- 338 [57] Justin N Wood and Samantha MW Wood. The development of newborn object recognition
339 in fast and slow visual worlds. *Proceedings of the Royal Society B: Biological Sciences*, 283
340 (1829):20160166, 2016.

- 341 [58] Samantha MW Wood, Scott P Johnson, and Justin N Wood. Automated study challenges the
342 existence of a foundational statistical-learning ability in newborn chicks. *Psychological Science*,
343 30(11):1592–1602, 2019.
- 344 [59] David D Cox, Philip Meier, Nadja Oertelt, and James J DiCarlo. ’breaking’ position-invariant
345 object recognition. *Nature Neuroscience*, 8(9):1145–1147, 2005.
- 346 [60] Taosheng Liu. Learning sequence of views of three-dimensional objects: The effect of temporal
347 coherence on object memory. *Perception*, 36(9):1320–1333, 2007.
- 348 [61] Guy Wallis, Benjamin T Backus, Michael Langer, Gesche Huebner, and Heinrich Bülthoff.
349 Learning illumination-and orientation-invariant representations of objects through temporal
350 association. *Journal of Vision*, 9(7):6–6, 2009.
- 351 [62] Nuo Li and James J DiCarlo. Unsupervised natural experience rapidly alters invariant object
352 representation in visual cortex. *Science*, 321(5895):1502–1507, 2008.
- 353 [63] Travis Meyer and Carl R Olson. Statistical learning of visual transitions in monkey infer-
354 otemporal cortex. *Proceedings of the National Academy of Sciences*, 108(48):19401–19406,
355 2011.
- 356 [64] Yasushi Miyashita. Neuronal correlate of visual associative long-term memory in the primate
357 temporal cortex. *Nature*, 335(6193):817–820, 1988.
- 358 [65] Giulio Matteucci and Davide Zoccolan. Unsupervised experience with temporal continuity of
359 the visual environment is causally involved in the development of v1 complex cells. *Science*
360 *Advances*, 6(22):eaba3742, 2020.
- 361 [66] Justin N Wood. A smoothness constraint on the development of object recognition. *Cognition*,
362 153:140–145, 2016.
- 363 [67] Justin N Wood, Aditya Prasad, Jason G Goldman, and Samantha MW Wood. Enhanced learning
364 of natural visual sequences in newborn chicks. *Animal Cognition*, 19:835–845, 2016.
- 365 [68] Justin N Wood and Samantha MW Wood. The development of invariant object recognition
366 requires visual experience with temporally smooth objects. *Cognitive Science*, 42(4):1391–1406,
367 2018.
- 368 [69] Samantha MW Wood and Justin N Wood. One-shot object parsing in newborn chicks. *Journal*
369 *of Experimental Psychology: General*, 150(11):2408, 2021.
- 370 [70] Arthur Aubret, Markus Ernst, Céline Teulière, and Jochen Triesch. Time to augment self-
371 supervised visual representation learning. *arXiv preprint arXiv:2207.13492*, 2022.
- 372 [71] Felix Schneider, Xia Xu, Markus R Ernst, Zhengyang Yu, and Jochen Triesch. Contrastive
373 learning through time. In *SVRHM 2021 Workshop@ NeurIPS*, 2021.
- 374 [72] Gerald M Edelman. Neural darwinism: selection and reentrant signaling in higher brain function.
375 *Neuron*, 10(2):115–125, 1993.
- 376 [73] Olaf Sporns and Gerald M Edelman. Solving bernstein’s problem: A proposal for the develop-
377 ment of coordinated movement by selection. *Child Development*, 64(4):960–981, 1993.
- 378 [74] Andrew Shtulman. Why people do not understand evolution: an analysis of the cognitive
379 barriers to fully grasping the unity of life. *Skeptic (Altadena, CA)*, 16(3):41–46, 2011.
- 380 [75] Richard C Lewontin. The units of selection. *Annual review of ecology and systematics*, 1(1):
381 1–18, 1970.
- 382 [76] Stephen Jay Gould. Darwinism and the expansion of evolutionary theory. *Science*, 216(4544):
383 380–387, 1982.
- 384 [77] Uri Hasson, Samuel A Nastase, and Ariel Goldstein. Direct fit to nature: an evolutionary
385 perspective on biological and artificial neural networks. *Neuron*, 105(3):416–434, 2020.
- 386 [78] Donsuk Lee, Samantha MW Wood, and Justin N Wood. Development of collective behavior in
387 newborn artificial agents. *arXiv preprint arXiv:2111.03796*, 2021.

388 **NeurIPS Paper Checklist**

389 **1. Claims**

390 Question: Do the main claims made in the abstract and introduction accurately reflect the
391 paper's contributions and scope?

392 Answer: [\[Yes\]](#) ,

393 Justification: Our main claim is that ViTs that learn with contrastive learning through time
394 show the same pattern of successes and failures as newborn chick (when trained with the
395 same training environment and tested with the same tasks). This is supported, in particular,
396 by Figure 2.

397 Guidelines:

- 398 • The answer NA means that the abstract and introduction do not include the claims
399 made in the paper.
- 400 • The abstract and/or introduction should clearly state the claims made, including the
401 contributions made in the paper and important assumptions and limitations. A No or
402 NA answer to this question will not be perceived well by the reviewers.
- 403 • The claims made should match theoretical and experimental results, and reflect how
404 much the results can be expected to generalize to other settings.
- 405 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
406 are not attained by the paper.

407 **2. Limitations**

408 Question: Does the paper discuss the limitations of the work performed by the authors?

409 Answer: [\[Yes\]](#)

410 Justification: Yes, we discuss the limitations of our work in Appendix Section A.7.2
411 (Limitations).

412 Guidelines:

- 413 • The answer NA means that the paper has no limitation while the answer No means that
414 the paper has limitations, but those are not discussed in the paper.
- 415 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 416 • The paper should point out any strong assumptions and how robust the results are to
417 violations of these assumptions (e.g., independence assumptions, noiseless settings,
418 model well-specification, asymptotic approximations only holding locally). The authors
419 should reflect on how these assumptions might be violated in practice and what the
420 implications would be.
- 421 • The authors should reflect on the scope of the claims made, e.g., if the approach was
422 only tested on a few datasets or with a few runs. In general, empirical results often
423 depend on implicit assumptions, which should be articulated.
- 424 • The authors should reflect on the factors that influence the performance of the approach.
425 For example, a facial recognition algorithm may perform poorly when image resolution
426 is low or images are taken in low lighting. Or a speech-to-text system might not be
427 used reliably to provide closed captions for online lectures because it fails to handle
428 technical jargon.
- 429 • The authors should discuss the computational efficiency of the proposed algorithms
430 and how they scale with dataset size.
- 431 • If applicable, the authors should discuss possible limitations of their approach to
432 address problems of privacy and fairness.
- 433 • While the authors might fear that complete honesty about limitations might be used by
434 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
435 limitations that aren't acknowledged in the paper. The authors should use their best
436 judgment and recognize that individual actions in favor of transparency play an impor-
437 tant role in developing norms that preserve the integrity of the community. Reviewers
438 will be specifically instructed to not penalize honesty concerning limitations.

439 **3. Theory Assumptions and Proofs**

440 Question: For each theoretical result, does the paper provide the full set of assumptions and
441 a complete (and correct) proof?

442 Answer: [NA]

443 Justification: Our work does not include any theoretical proofs or mathematical derivations.

444 Guidelines:

- 445 • The answer NA means that the paper does not include theoretical results.
- 446 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
447 referenced.
- 448 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 449 • The proofs can either appear in the main paper or the supplemental material, but if
450 they appear in the supplemental material, the authors are encouraged to provide a short
451 proof sketch to provide intuition.
- 452 • Inversely, any informal proof provided in the core of the paper should be complemented
453 by formal proofs provided in appendix or supplemental material.
- 454 • Theorems and Lemmas that the proof relies upon should be properly referenced.

445 4. Experimental Result Reproducibility

456 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
457 perimental results of the paper to the extent that it affects the main claims and/or conclusions
458 of the paper (regardless of whether the code and data are provided or not)?

459 Answer: [Yes]

460 Justification: We provide detailed instructions for each experiment in Sections 3.1, 3.2, and
461 3.3, and instructions for generating the datasets in Appendix A.2. Additionally, information
462 on training and testing the models can be found in Appendix Sections A.3, A.4, A.5, and
463 A.6. Code and data for reproducibility are provided in Appendix Section A.8.

464 Guidelines:

- 465 • The answer NA means that the paper does not include experiments.
- 466 • If the paper includes experiments, a No answer to this question will not be perceived
467 well by the reviewers: Making the paper reproducible is important, regardless of
468 whether the code and data are provided or not.
- 469 • If the contribution is a dataset and/or model, the authors should describe the steps taken
470 to make their results reproducible or verifiable.
- 471 • Depending on the contribution, reproducibility can be accomplished in various ways.
472 For example, if the contribution is a novel architecture, describing the architecture fully
473 might suffice, or if the contribution is a specific model and empirical evaluation, it may
474 be necessary to either make it possible for others to replicate the model with the same
475 dataset, or provide access to the model. In general, releasing code and data is often
476 one good way to accomplish this, but reproducibility can also be provided via detailed
477 instructions for how to replicate the results, access to a hosted model (e.g., in the case
478 of a large language model), releasing of a model checkpoint, or other means that are
479 appropriate to the research performed.
- 480 • While NeurIPS does not require releasing code, the conference does require all submis-
481 sions to provide some reasonable avenue for reproducibility, which may depend on the
482 nature of the contribution. For example
 - 483 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
484 to reproduce that algorithm.
 - 485 (b) If the contribution is primarily a new model architecture, the paper should describe
486 the architecture clearly and fully.
 - 487 (c) If the contribution is a new model (e.g., a large language model), then there should
488 either be a way to access this model for reproducing the results or a way to reproduce
489 the model (e.g., with an open-source dataset or instructions for how to construct
490 the dataset).
 - 491 (d) We recognize that reproducibility may be tricky in some cases, in which case
492 authors are welcome to describe the particular way they provide for reproducibility.
493 In the case of closed-source models, it may be that access to the model is limited in

494 some way (e.g., to registered users), but it should be possible for other researchers
495 to have some path to reproducing or verifying the results.

496 5. Open access to data and code

497 Question: Does the paper provide open access to the data and code, with sufficient instruc-
498 tions to faithfully reproduce the main experimental results, as described in supplemental
499 material?

500 Answer: [Yes]

501 Justification: We provide open access to our datasets and the models with detailed instruc-
502 tions on our GitHub page in Appendix Section A.8.

503 Guidelines:

- 504 • The answer NA means that paper does not include experiments requiring code.
- 505 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
506 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 507 • While we encourage the release of code and data, we understand that this might not be
508 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
509 including code, unless this is central to the contribution (e.g., for a new open-source
510 benchmark).
- 511 • The instructions should contain the exact command and environment needed to run to
512 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
513 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 514 • The authors should provide instructions on data access and preparation, including how
515 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 516 • The authors should provide scripts to reproduce all experimental results for the new
517 proposed method and baselines. If only a subset of experiments are reproducible, they
518 should state which ones are omitted from the script and why.
- 519 • At submission time, to preserve anonymity, the authors should release anonymized
520 versions (if applicable).
- 521 • Providing as much information as possible in supplemental material (appended to the
522 paper) is recommended, but including URLs to data and code is permitted.

523 6. Experimental Setting/Details

524 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
525 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
526 results?

527 Answer: [Yes]

528 Justification: We provide all experimental details relevant to replicating our results in
529 the Appendix. Additionally, we also provide information on our computational models
530 (architecture and objective function) in Section 2 (Methods).

531 Guidelines:

- 532 • The answer NA means that the paper does not include experiments.
- 533 • The experimental setting should be presented in the core of the paper to a level of detail
534 that is necessary to appreciate the results and make sense of them.
- 535 • The full details can be provided either with the code, in appendix, or as supplemental
536 material.

537 7. Experiment Statistical Significance

538 Question: Does the paper report error bars suitably and correctly defined or other appropriate
539 information about the statistical significance of the experiments?

540 Answer: [Yes]

541 Justification: We include error bars (showing Standard Error) in all of our bar charts (reported
542 in Fig 1, Panel 3).

543 Guidelines:

- 544 • The answer NA means that the paper does not include experiments.

- 545 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
546 dence intervals, or statistical significance tests, at least for the experiments that support
547 the main claims of the paper.
- 548 • The factors of variability that the error bars are capturing should be clearly stated (for
549 example, train/test split, initialization, random drawing of some parameter, or overall
550 run with given experimental conditions).
- 551 • The method for calculating the error bars should be explained (closed form formula,
552 call to a library function, bootstrap, etc.)
- 553 • The assumptions made should be given (e.g., Normally distributed errors).
- 554 • It should be clear whether the error bar is the standard deviation or the standard error
555 of the mean.
- 556 • It is OK to report 1-sigma error bars, but one should state it. The authors should
557 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
558 of Normality of errors is not verified.
- 559 • For asymmetric distributions, the authors should be careful not to show in tables or
560 figures symmetric error bars that would yield results that are out of range (e.g. negative
561 error rates).
- 562 • If error bars are reported in tables or plots, The authors should explain in the text how
563 they were calculated and reference the corresponding figures or tables in the text.

564 8. Experiments Compute Resources

565 Question: For each experiment, does the paper provide sufficient information on the com-
566 puter resources (type of compute workers, memory, time of execution) needed to reproduce
567 the experiments?

568 Answer: [Yes]

569 Justification: We provide all the relevant information in Appendix Section A.6 (Training
570 Details).

571 Guidelines:

- 572 • The answer NA means that the paper does not include experiments.
- 573 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
574 or cloud provider, including relevant memory and storage.
- 575 • The paper should provide the amount of compute required for each of the individual
576 experimental runs as well as estimate the total compute.
- 577 • The paper should disclose whether the full research project required more compute
578 than the experiments reported in the paper (e.g., preliminary or failed experiments that
579 didn't make it into the paper).

580 9. Code Of Ethics

581 Question: Does the research conducted in the paper conform, in every respect, with the
582 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

583 Answer: [Yes]

584 Justification: Yes, all author(s) in this paper have reviewed the NeurIPS Code of Ethics and
585 abide by its rules.

586 Guidelines:

- 587 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 588 • If the authors answer No, they should explain the special circumstances that require a
589 deviation from the Code of Ethics.
- 590 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
591 eration due to laws or regulations in their jurisdiction).

592 10. Broader Impacts

593 Question: Does the paper discuss both potential positive societal impacts and negative
594 societal impacts of the work performed?

595 Answer: [Yes]

596 Justification: We include a Broader Impacts section, but our societal impacts are limited
597 because we are reporting foundational research that is not tied to particular applications, let
598 alone deployments.

599 Guidelines:

- 600 • The answer NA means that there is no societal impact of the work performed.
- 601 • If the authors answer NA or No, they should explain why their work has no societal
602 impact or why the paper does not address societal impact.
- 603 • Examples of negative societal impacts include potential malicious or unintended uses
604 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
605 (e.g., deployment of technologies that could make decisions that unfairly impact specific
606 groups), privacy considerations, and security considerations.
- 607 • The conference expects that many papers will be foundational research and not tied
608 to particular applications, let alone deployments. However, if there is a direct path to
609 any negative applications, the authors should point it out. For example, it is legitimate
610 to point out that an improvement in the quality of generative models could be used to
611 generate deepfakes for disinformation. On the other hand, it is not needed to point out
612 that a generic algorithm for optimizing neural networks could enable people to train
613 models that generate Deepfakes faster.
- 614 • The authors should consider possible harms that could arise when the technology is
615 being used as intended and functioning correctly, harms that could arise when the
616 technology is being used as intended but gives incorrect results, and harms following
617 from (intentional or unintentional) misuse of the technology.
- 618 • If there are negative societal impacts, the authors could also discuss possible mitigation
619 strategies (e.g., gated release of models, providing defenses in addition to attacks,
620 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
621 feedback over time, improving the efficiency and accessibility of ML).

622 11. Safeguards

623 Question: Does the paper describe safeguards that have been put in place for responsible
624 release of data or models that have a high risk for misuse (e.g., pretrained language models,
625 image generators, or scraped datasets)?

626 Answer: [NA]

627 Justification: The datasets and models used in this study pose no risks. The datasets consist
628 of virtual 3D and 2D objects that do not pose any safety concerns.

629 Guidelines:

- 630 • The answer NA means that the paper poses no such risks.
- 631 • Released models that have a high risk for misuse or dual-use should be released with
632 necessary safeguards to allow for controlled use of the model, for example by requiring
633 that users adhere to usage guidelines or restrictions to access the model or implementing
634 safety filters.
- 635 • Datasets that have been scraped from the Internet could pose safety risks. The authors
636 should describe how they avoided releasing unsafe images.
- 637 • We recognize that providing effective safeguards is challenging, and many papers do
638 not require this, but we encourage authors to take this into account and make a best
639 faith effort.

640 12. Licenses for existing assets

641 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
642 the paper, properly credited and are the license and terms of use explicitly mentioned and
643 properly respected?

644 Answer: [Yes]

645 Justification: Yes, all the original assets used in this study (if any) are properly cited.

646 Guidelines:

- 647 • The answer NA means that the paper does not use existing assets.
- 648 • The authors should cite the original paper that produced the code package or dataset.

- 649 • The authors should state which version of the asset is used and, if possible, include a
650 URL.
- 651 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 652 • For scraped data from a particular source (e.g., website), the copyright and terms of
653 service of that source should be provided.
- 654 • If assets are released, the license, copyright information, and terms of use in the
655 package should be provided. For popular datasets, `paperswithcode.com/datasets`
656 has curated licenses for some datasets. Their licensing guide can help determine the
657 license of a dataset.
- 658 • For existing datasets that are re-packaged, both the original license and the license of
659 the derived asset (if it has changed) should be provided.
- 660 • If this information is not available online, the authors are encouraged to reach out to
661 the asset's creators.

662 13. **New Assets**

663 Question: Are new assets introduced in the paper well documented and is the documentation
664 provided alongside the assets?

665 Answer: [Yes]

666 Justification: We provide open access to all the assets used in this paper (models, datasets,
667 and scripts) in Appendix Section A.8.

668 Guidelines:

- 669 • The answer NA means that the paper does not release new assets.
- 670 • Researchers should communicate the details of the dataset/code/model as part of their
671 submissions via structured templates. This includes details about training, license,
672 limitations, etc.
- 673 • The paper should discuss whether and how consent was obtained from people whose
674 asset is used.
- 675 • At submission time, remember to anonymize your assets (if applicable). You can either
676 create an anonymized URL or include an anonymized zip file.

677 14. **Crowdsourcing and Research with Human Subjects**

678 Question: For crowdsourcing experiments and research with human subjects, does the paper
679 include the full text of instructions given to participants and screenshots, if applicable, as
680 well as details about compensation (if any)?

681 Answer: [NA]

682 Justification: Our study does not involve research with Human Subjects.

683 Guidelines:

- 684 • The answer NA means that the paper does not involve crowdsourcing nor research with
685 human subjects.
- 686 • Including this information in the supplemental material is fine, but if the main contribu-
687 tion of the paper involves human subjects, then as much detail as possible should be
688 included in the main paper.
- 689 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
690 or other labor should be paid at least the minimum wage in the country of the data
691 collector.

692 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 693 Subjects**

694 Question: Does the paper describe potential risks incurred by study participants, whether
695 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
696 approvals (or an equivalent approval/review based on the requirements of your country or
697 institution) were obtained?

698 Answer: [NA]

699 Justification: Our study does not involve research with Human Subjects. Our study does not
700 require IRB Approvals or Equivalent.

701
702
703
704
705
706
707
708
709
710
711

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

712 **A Appendix**

713 **A.1 Experiment 4: Comparing Different Objective Functions**

714 In Experiments 1-3, we compared CNNs and ViTs that had the same temporal learning objective.
 715 We used contrastive learning through time because it implements the UTL principle discovered in
 716 neuroscience and behavioral experiments. In Experiment 4, we assessed the contribution of the
 717 temporal learning objective by comparing CNNs and ViTs across other objective functions. If the
 718 CNNs and ViTs still show the same pattern of performance (i.e., ViTs, but not CNNs, are impaired
 719 when learning from line drawings), then the architecture alone would be the main contributing factor
 720 for mimicking visual learning in chicks. However, if the pattern changes, then both the architecture
 721 and the objective function would be essential for mimicking learning in chicks.

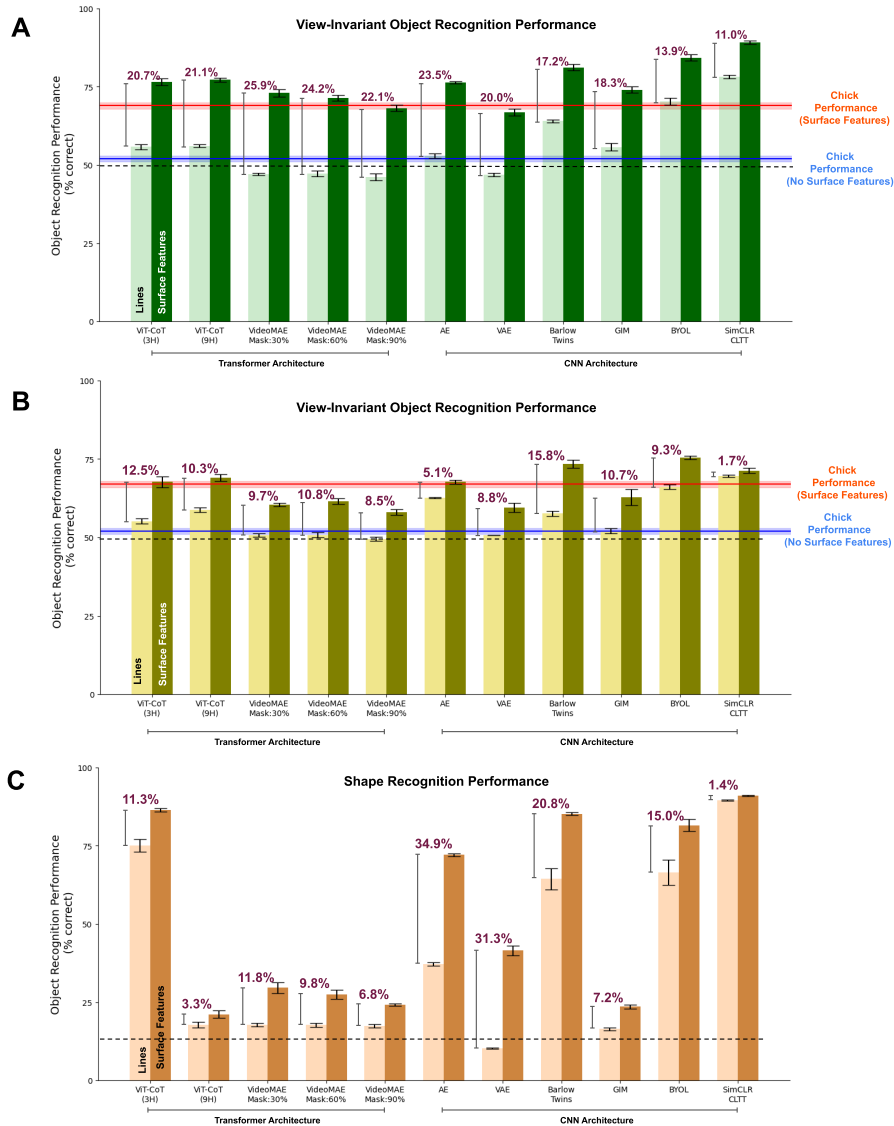


Figure S 1: Object recognition performance of a range of CNNs and ViTs in (A) Experiment 1, (B) Experiment 2, and (C) Experiment 3.

722 We repeated Experiments 1-3 with four additional ViT models and five additional CNN models. All
 723 of the models used different self-supervised objective functions. As shown in Fig. S1, most of the

724 ViT and CNN models were significantly impaired when learning from line drawings compared to
 725 learning from objects with surface features. Yet, many of the CNN models still succeeded in both
 726 conditions, learning object recognition even from line drawings (unlike chicks).

727 Overall, Experiment 4 shows that researchers will need to consider both the architecture and the
 728 objective function to build models of visual learning. We show that, by precisely characterizing both
 729 the architecture (transformer) *and* the objective function (temporal learning), deep neural networks
 730 can serve as accurate image-computable models of visual learning.

731 **A.2 Data Generation**

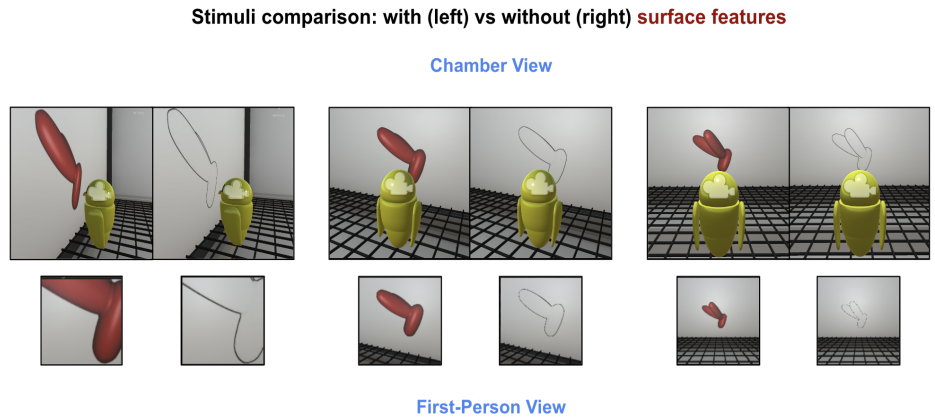


Figure S 2: The virtual chamber in the normal object and line drawing conditions. (Top) The agent visually explores the chamber, randomly moving from place to another. (Bottom) First-person images captured from the camera attached to the agent's head. We use the first-person images to train the ViTs and CNNs.

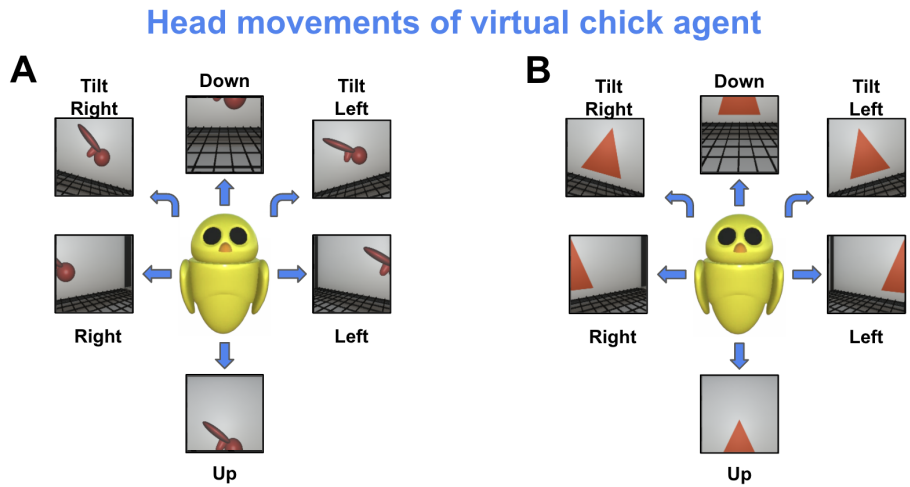


Figure S 3: Head movements of the virtual agent across the three axes of rotations (yaw, roll, and tilt). The agent moved its head 60° on each axis. The images show how head movements provide a form of natural data augmentation.

732 We created a virtual animal chamber in the Unity Game Engine (Fig. S2). The virtual chamber had
 733 two 19" LCD monitors on opposite sides, while the other two sides of the chamber were white walls.
 734 The LCDs were used to display virtual 3D or 2D objects moving on a white background at the center
 735 of the screens. The floor of the virtual chamber was constructed with black wire mesh and had a
 736 provision for food and water next to one of the chamber walls. The dimensions of the virtual chamber

737 were 66 cm (L) x 42 cm (W) x 69 cm (H). The chamber also contained a virtual chick agent; the
 738 dimensions of the chick agent were 3.5 units (H) x 1.2 units (L).

739 This virtual chamber was equipped with two cameras. One camera was placed in the position of
 740 the agent’s eyes to capture first-person RGB images, simulating the visual experiences of newborn
 741 chicks. The second camera was placed on the chamber’s ceiling to capture a top view of the agent’s
 742 movement. To simulate the visual diet available in the chambers, the agent moved to random locations
 743 inside the chamber, at a speed of 1.5 units per second. While moving, the agent maintained a constant
 744 gaze at the object. Once it reached its destination, the agent then moved its head along all three axes
 745 (yaw, roll, and tilt) in a random order (Fig. S3). These head movements lasted for 9.5 seconds. This
 746 cycle was repeated until 80,000 first-person images were collected in each rearing condition. The
 747 same method was used to collect test data, except the agent kept their gaze fixed on the object.

748 This simulation approach canvassed the range of visual experiences that chicks could acquire in the
 749 chamber. The approach did not directly simulate a specific chick’s visual experiences. The approach
 750 also did not capture views chicks may have seen of their own bodies (e.g., wings, feet). Our virtual
 751 agent could not see its body, so its visual diet was limited to views of the chamber. As such, this
 752 approach established a baseline of what could be learned when a model has access to the same visual
 753 environment as newborn chicks.

754 A.3 Architectures

755 We report all the model architectures and their hyperparameters in Table 1.

Table 1: Architectures and Hyperparameters for various self-supervised learning models

Model	Parameters (M)	Attention Heads	Layers	Learning Objective	Batch Size
ViT-3H	16.9	3	3	Contrastive Learning Through Time	128
ViT-9H	59.4	9	9	Contrastive Learning Through Time	128
VideoMAE-0.3	53.9	6	6	Video Reconstruction	32x8(GPUs)
VideoMAE-0.6	53.9	6	6	Video Reconstruction	32x8(GPUs)
VideoMAE-0.9	53.9	6	6	Video Reconstruction	32x8(GPUs)
SimCLR-CLTT	7.9	NA	10	Contrastive Learning Through Time	512
BYOL	15.9	NA	10	Asymmetric Embedding	512
Barlow Twins	7.9	NA	10	Joint Embedding	512
AE	15.5	NA	10	Image Reconstruction	128
VAE	15.6	NA	10	Image Reconstruction	128
GIM	16.5	NA	10	Non-Backpropagation Contrastive Learning	32

756 A.3.1 ViT-CoT

757 We systematically varied the number of attention heads and transformer layers to create different
 758 architecture sizes for ViTs. For instance, we used three attention heads and layers to create ViT-3H.
 759 For ViT-9H, we increased the number of attention heads and layers to nine. The last layer of the
 760 ViT-CoTs generated a 512-dimensional embedding, which was then passed through the loss function.
 761 Each architecture was trained using self-supervised learning with a contrastive learning through
 762 time objective function. To preserve the temporal relationships between consecutive frames, we did
 763 not shuffle the frames in the dataset. Additionally, to avoid hardcoding spatial knowledge in the

764 ViT-CoTs, we did not use any convolutional layers to generate image patches. The models were
765 trained using images of size 64x64 and a patch size of 8x8. A constant learning rate of 0.0001 was
766 used to train the models.

767 A.3.2 VideoMAE

768 In the VideoMAE architecture, both the encoder and decoder blocks had six layers and attention heads.
769 The VideoMAEs were trained by sampling 16 frames from the training set with a temporal stride of 1.
770 Each batch sample had dimensions of (16x3x64x64), where 16 represents the temporal window and
771 3x64x64 indicates the image dimensions. Subsequently, a random mask of spatial dimension 8x8 and
772 temporal dimension of 2 (2x8x8) was applied to the training batch. The visible patches (non-masked
773 patches) were encoded by the VideoMAE encoder and passed on to the VideoMAE decoder. The
774 decoder combined the encoded features and the masked patches to reconstruct the entire sequence of
775 temporal frames. We experimented with three masking ratios: 30%, 60%, and 90%.

776 A.3.3 CNN

777 For the CNN models, we created a custom ResNet architecture (ResNet-10). Each architecture
778 consisted of two residual blocks, totaling of 10 convolutional layers. We used the same bridge
779 connections between the residual blocks as implemented in default ResNets. Similar to the ViT-CoTs,
780 the last layer of the CNNs generated a 512-dimensional embedding, which was then passed through
781 the loss function. To train SimCLR-CLTT, we used a learning rate scheduler with the warm-up
782 epochs set to 5. Additionally, to preserve the temporal relationships between consecutive frames, we
783 did not shuffle the frames in the dataset.

784 A.4 Objective Function

785 Many behavioral studies provide evidence that human adults use UTL to learn object representations
786 [59–61]. UTL has also been found on the neurophysiological level in adult monkeys [62–64].
787 There is even evidence that newborn animals (including chicks) use UTL to build their first object
788 representation [65, 66, 57, 67–69]. These findings suggest that UTL is foundational to visual learning.

789 To incorporate UTL in our models, we used a temporal learning algorithm, Contrastive Learning
790 Through Time (CLTT), that can be implemented in both CNNs [70, 71] and ViTs [8]. CLTT leverages
791 the temporal structure of natural visual experience, without relying on supervision or labeled data
792 (see Fig. S4). The algorithm contrasts temporally adjacent instances (positive examples) against non-
793 adjacent instances (negative examples), thereby learning representations that capture the underlying
794 dynamics, context, and patterns across time.

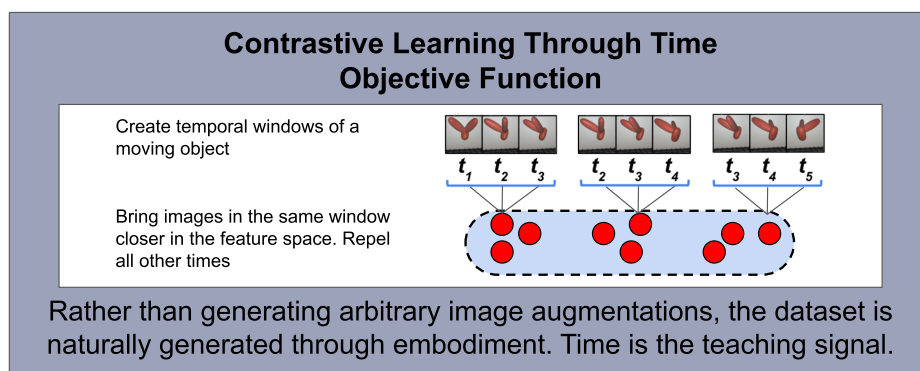


Figure S 4: Contrastive Learning Through Time (CLTT) objective function used with the SimCLR-CLTT (CNN) and ViT-COT (transformer) models. The algorithm pushes together features that occur in the same temporal window (300 ms time window), akin to the 100-400 ms spike-timing-dependent plasticity temporal learning window in brains.

795 **A.5 Evaluation**

796 After training the models (encoders), we evaluated their classification performance using the
797 stimuli. Task performance was assessed by removing the last fully connected layer of the network,
798 adding a new fully connected linear readout layer on top of the last layer of each trained encoder, and
799 then training only the parameters of the readout layer on the object classification task. The linear
800 classifiers contained 512 neurons, each of which received input from one of the 512 neurons in the
801 final layer of the model. The linear readout layers were optimized for binary cross-entropy loss.

802 To train and test the linear classifiers, we used the test images collected from the agents moving
803 through the virtual chambers (10,000 images for each of two objects across 12 viewpoint ranges,
804 see Fig. S5). When training the linear classifiers, the object identities were used as the ground-truth
805 labels. Since the encoder weights were frozen, the supervised training of the linear classifiers did not
806 change the features learned by the model.

807 To evaluate whether the features learned by the models could generalize across novel viewpoints, we
808 used a cross-validated K-fold analysis to train/test the linear classifiers, where each fold contained
809 images from one of the 12 viewpoint ranges. Specifically, the test images were divided into 12 folds,
810 with each fold containing images of each object rotating through 1 viewpoint range. The linear
811 classifiers were cross-validated by training on 11 folds (11 viewpoint ranges) and testing on the
812 held-out fold (1 viewpoint range).

813 The linear classifiers were trained on 11,000 total images. During training, we used a batch size
814 of 128 for 100 epochs. Transfer performance was evaluated by first fitting the parameters of the
815 linear classifier on the training set and then measuring classification accuracy on the held-out test set.
816 We report average cross-validated performance on the held-out images not used to train the linear
817 readout layer. Thus, all of our results reflect the generalization performance of the models across
818 novel viewpoints.

819 In Experiment 2, we reused the same linear classifier design from Experiment 1 to conduct binary
820 classification between the two objects. In both experiments, the linear classifier was always trained
821 on 10,000 samples.

822 In Experiment 3, the linear classifier had 8 output neurons, each corresponding to an object class.
823 We used softmax and categorical cross-entropy loss to train the linear classifier. The training dataset
824 consisted of 8 object classes with each class having 2,500 samples. To construct a test set, we split
825 the training set in half by selecting the initial 1,250 samples from each class. This way, the linear
classifier could be trained and evaluated on 10,000 samples (1250 samples x 8 classes).

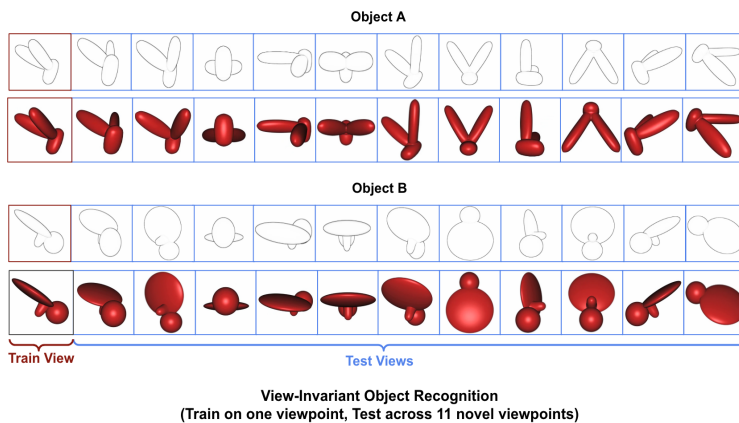
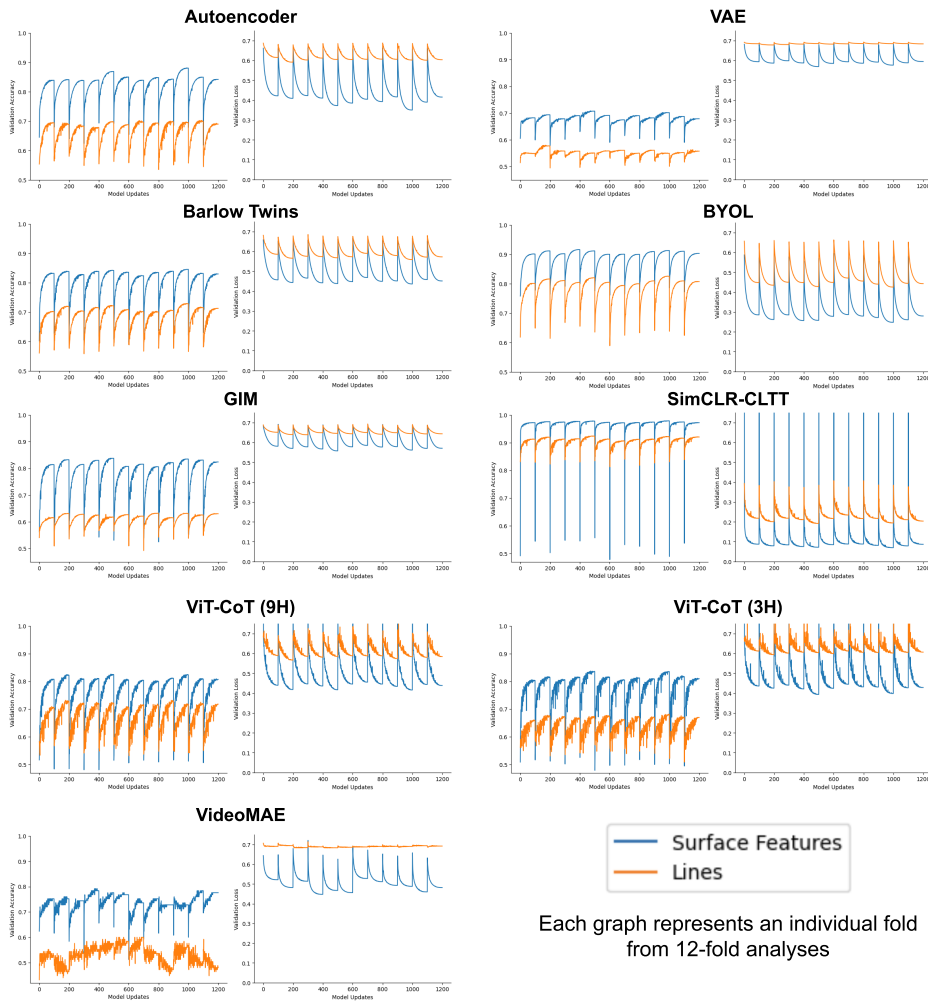


Figure S 5: Viewpoints used in Experiment 1 for the view-invariant object recognition task. The encoder was trained on a single viewpoint and tested on 11 novel viewpoints, using a 12-fold cross-validation design with a linear classifier. The images show object images for the line drawings (top) and normal objects with surface features (bottom).

Validation Accuracy and Validation Loss of Linear Probe Feature Extraction for a K-Fold Analysis



A linear probe effectively extracts representations from an encoder trained on **rich surface features** compared to the sparse representation obtained from the encoder trained on line drawings.

Figure S 6: Linear classifier evaluation results for Experiment 4. The plots show the validation accuracy and validation loss for the linear classifiers attached to nine different visual encoders. Normal objects with surface features consistently provide stronger and more robust learning signals than line drawings.

827 A.5.1 Evaluation Results

828 In Fig. S6, we present the validation accuracy and validation loss data for the linear classifiers
 829 trained on frozen encoders in Experiment 1. We observed a consistent pattern across all models: the
 830 validation accuracy was high when the linear classifier was evaluated on a frozen encoder trained on
 831 normal objects, but it was low when evaluated on line drawings. Similarly, the validation loss
 832 was low for encoders trained on normal objects compared to those trained on line drawings. This
 833 indicates that linear classifiers can effectively extract rich surface features, but struggle to disentangle
 834 features when the encoders are trained on line drawings.

835 Fig. S7 compares the validation accuracy and validation loss for ViT and SimCLR, both having the
 836 same objective function (contrastive learning through time) but different architectures in case

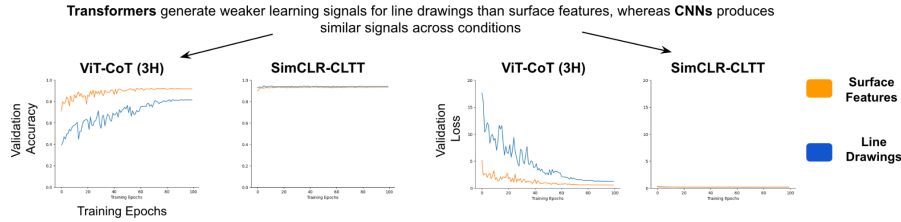


Figure S 7: Linear classifier evaluation results for Experiment 3. The plots compare ViTs and CNNs that have the same contrastive learning through time objective function. Transformers generate weaker learning signals when reared with line drawings versus normal objects (like chicks), whereas CNNs produce similar learning signals for normal objects and line drawings.

837 of experiment 3. The transformers provide strong training signals to the linear classifier when the
 838 encoder is trained on surface features, as opposed to line drawings. In contrast, the CNN backbone
 839 provides consistent training signals regardless of whether it is trained on surface features or line
 840 drawings.

841 A.6 Training Details

842 We trained each model using 3 different seeds and 100 epochs. All models, except VideoMAEs,
 843 were trained on a single NVIDIA A10 GPU. VideoMAEs were trained using multi-GPU distributed
 844 training across 8 NVIDIA A10 GPUs. Each GPU had 24 gigabytes of memory. We report the number
 845 of trainable parameters for each model in Table 1.

846 A.7 Discussion

847 Our study provides a new form of guidance for building brain-like AI systems. Researchers have long
 848 attempted to build machines that learn like brains, but almost all prior studies compared animals and
 849 machines that were raised (trained) in different environments. If animals and machines learn from
 850 different training data, then it is impossible to determine whether machines learn like brains (i.e.,
 851 differences in performance could be due to the algorithm, training data, or some combination of the
 852 factors). We overcome this barrier by performing parallel controlled-rearing experiments on newborn
 853 chicks and AI algorithms, matching training data across animals and machines. This allowed us to
 854 distinguish between candidate model classes (ViTs vs. CNNs) and discover AI systems (ViTs) that
 855 show the same learning outcomes as newborn brains. We found that transformers, which are typically
 856 considered to be less "brain-like" than CNNs, are the more accurate model of visual learning.

857 A.7.1 Theoretical simulations of the origins of vision

858 There is a long history of attempts to characterize the core learning machinery underlying intelligence.
 859 Our work expands earlier techniques [72, 73] that used theoretical simulations to study whether
 860 blind, evolution-like fitting processes can explain the rapid, self-organized development of visual
 861 intelligence. Earlier simulations were limited by compute power, so researchers could not run
 862 *image computable* simulations testing whether core visual skills really are learnable via generic
 863 fitting machinery. Image computable simulations are essential for testing fitting theories of brain
 864 development because the outcomes of evolutionary processes can be counter intuitive [74].

865 Transformers are ideal models for running simulations of evolution-like fitting processes. The ma-
 866 chinery underlying evolution involves blind, brute-force fitting, in which a generic high-dimensional
 867 combinatorial medium (the genetic code) adapts to the environment [75, 76]. Likewise, transformers
 868 are blind, brute-force fitting systems, in which a generic high-dimensional combinatorial medium
 869 (neural networks) adapts to the data distributions in the environment. Both natural selection and
 870 transformers start from scratch and produce complex animal forms (natural selection) and mental
 871 skills (transformers). Since evolution and transformers operate by common general fitting principles
 872 [77], they can be united under a common framework.

873 We show that transformers, which start from scratch (no prior knowledge of objects or space) and
 874 learn through blind fitting, are sufficient to account for successes and failures of visual object learning

875 in newborn chicks. Based on these (and other [8, 78, 77]) findings, we speculate that learning in
876 the brain can be understood in evolutionary terms, as a dynamic high-dimensional system adapting
877 (fitting) to the spatiotemporal data distributions underlying sensory experiences. Under this view,
878 object recognition is not a hard-coded system, structure, primitive, module, or program; rather, it
879 is an emergent property of generic temporal fitting machinery adapting to the embodied visual data
880 streams acquired by newborn animals.

881 **A.7.2 Limitations**

882 One limitation of our study is the models were trained passively, learning from batches of images in a
883 pre-specified order. Newborn animals, in contrast, interact with their environment to produce their
884 own training data. Future studies could close this gap between animals and machines by embodying
885 CNNs and ViTs in artificial agents that collect their own training data from the environment. A
886 second limitation is we do not know *why* the objects with surface features provide better learning
887 signals than line drawings. In Appendix Section A.5.1, we provide preliminary results showing that
888 objects with surface features provide more robust learning signals than line drawings.

889 **A.7.3 Broader Impact**

890 This paper tackles a question at the heart of cognitive science: What are the core learning algorithms
891 in brains? By demonstrating that transformers produce similar learning outcomes as newborn animals,
892 our work shows that transformers can be powerful modeling tools for studying how brains learn
893 to perceive and understand the world. Our work also provides an important step towards building
894 “naturally intelligent” learning systems. Naturally intelligent learning algorithms are an untapped
895 goldmine for inspiring the next generation of machine learning systems.

896 **A.8 Data and Code Availability**

897 The code and data needed to reproduce these findings will be available upon publication.