

---

# Causality $\neq$ Decodability, and Vice Versa: Lessons from Interpreting Counting ViTs

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Mechanistic interpretability seeks to uncover how internal components of neural  
2 networks give rise to predictions. A persistent challenge, however, is disentangling  
3 two often conflated notions: decodability—the recoverability of information from  
4 hidden states—and causality—the extent to which those states functionally influ-  
5 ence outputs. In this work, we investigate their relationship in vision transformers  
6 (ViTs) fine-tuned for object counting. Using activation patching, we test the causal  
7 role of spatial and CLS tokens by transplanting activations across clean–corrupted  
8 image pairs. In parallel, we train linear probes to assess the decodability of count  
9 information at different depths. Our results reveal systematic mismatches: middle-  
10 layer object tokens exert strong causal influence despite being weakly decodable,  
11 whereas final-layer object tokens support accurate decoding yet are functionally  
12 inert. Similarly, the CLS token becomes decodable in mid-layers but only ac-  
13 quires causal power in the final layers. These findings highlight that decodability  
14 and causality reflect complementary dimensions of representation—what informa-  
15 tion is present versus what is used—and that their divergence can expose hidden  
16 computational circuits.

## 17 1 Introduction

18 Mechanistic interpretability seeks to uncover how internal components of neural networks contribute  
19 to predictions, moving beyond aggregate performance metrics toward causal understanding of model  
20 behavior [2, 17, 4, 5, 15, 21, 9, 13, 14]. A central challenge is distinguishing between two notions  
21 often conflated in practice: the decodability of information from representations, and the causal  
22 use of that information by the model. Decodability methods, such as linear probing, test whether  
23 a variable of interest can be recovered from hidden states [1]. Causal methods, such as activation  
24 patching, instead test whether modifying activations changes the model’s outputs [6, 22]. While both  
25 approaches provide valuable perspectives, it remains unclear how they align—or diverge—across  
26 layers and token types in large models.

27 This distinction is especially pertinent for vision transformers (ViTs), whose predictions arise from the  
28 interaction of local patch embeddings and a global classification (CLS) token. In ViTs, local patches  
29 may contain object-specific information, while the CLS token aggregates global scene evidence. Yet  
30 whether information contained in these representations is actually used in making predictions is less  
31 well understood. For example, a token may carry highly decodable features but exert no influence  
32 on the output, or conversely, a token with weakly decodable information may nonetheless causally  
33 drive predictions when perturbed. Understanding this gap is crucial for accurately characterizing  
34 what ViTs represent and how they compute.

35 In this work, we investigate the relationship between decodability and causality in a vision transformer  
36 fine-tuned for object counting [3, 11, 12]. Using activation patching, we transplant hidden activations

37 from clean and corrupted image pairs to test which tokens influence predictions at different depths. In  
 38 parallel, we train linear probes on object patches, CLS tokens, and background patches to assess their  
 39 decodability. Comparing the two perspectives reveals systematic mismatches. Middle-layer object  
 40 tokens, though weakly decodable, exert strong causal influence when patched. By contrast, final-layer  
 41 object tokens support highly accurate decoding yet are functionally inert, with predictions unaffected  
 42 by patching. Similarly, CLS tokens become decodable in the middle layers but only acquire causal  
 43 power in the final layers.

44 Our findings highlight that decodability and causality are not interchangeable lenses on model  
 45 behavior. Instead, they reflect complementary dimensions of representation: what information is  
 46 present, and what information is used. By demonstrating their divergence in a concrete setting, we  
 47 argue that both are necessary for a comprehensive interpretability analysis, and that mismatches  
 48 between them may reveal hidden computational circuits.

## 49 2 Activation Patching for Causality

50 Activation patching [6, 22] is a causal interpretability method that tests whether specific activations  
 51 are used by the model to make predictions. Given a source input  $x_s$  and a target input  $x_t$ , we forward  
 52 each through the model to obtain hidden states  $h_s^l, h_t^l$  at layer  $l$ . We then form a patched run by  
 53 replacing part of the target activations with those from the source:

$$\tilde{h}_t^l = h_t^l \text{ with } h_{t,i}^l \leftarrow h_{s,i}^l,$$

54 where  $i$  indexes the token or component being patched. Continuing the forward pass from  $\tilde{h}_t^l$  yields a  
 55 new output  $\tilde{f}(x_t)$ . If  $\tilde{f}(x_t)$  shifts toward  $f(x_s)$ , then the patched activations causally influence the  
 56 prediction.

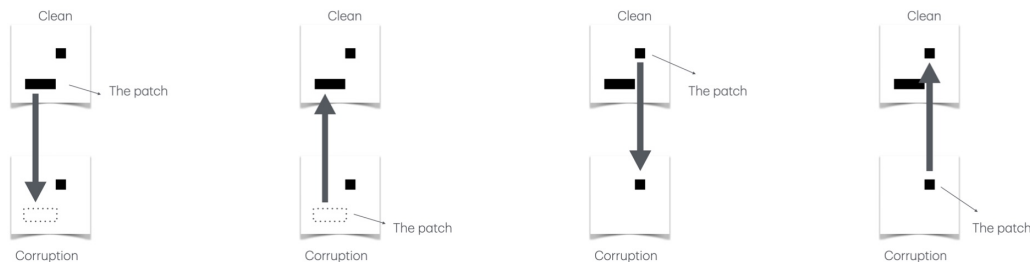


Figure 1: Activation patching experiments, corresponding to experiments 1 through 4, where object A is the  $1 \times 1$  square and object B is the  $1 \times 3$  rectangle.

57 As an illustrative example of our activation patching setup, we consider an image containing two  
 58 objects, A and B (the clean image). We construct a corrupted counterpart by removing object B.  
 59 When passed through the counting-finetuned vision transformer, the clean image yields a prediction  
 60 of 2, while the corrupted image yields a prediction of 1, as expected (see Appendix B for details on  
 61 the model and dataset). We then perform the following activation patching experiments (Figure 1):

- 62 1. Patch the activation of object B from the clean run into the corrupted run.
- 63 2. Patch the activation of the empty patches in the corrupted run (corresponding to the patch  
 64 locations of object B in the clean run) into the clean run.
- 65 3. Patch the activation of object A from the clean run into the corrupted run.
- 66 4. Patch the activation of object A from the corrupted run into the clean run.
- 67 5. Patch the CLS token from the clean run into the corrupted run.
- 68 6. Patch the CLS token from the corrupted run into the clean run.

69 In our setup, activations are patched layer-wise: an activation from a given layer in the source run is  
 70 transplanted into the corresponding layer of the target run. To assess how information at different

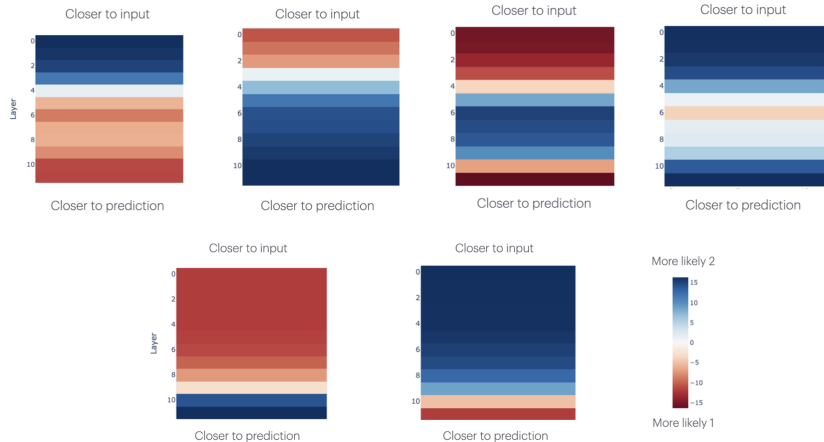


Figure 2: Activation patching results, corresponding to experiments 1 through 6 in left-right, top-bottom order.

71 depths influences model behavior, we measure the effect of patching by computing the logit difference  
 72 between the model’s predictions for class 2 versus class 1 as a function of the patched layer. The  
 73 results are shown in Figure 2.

74 In experiment 1, patching a token corresponding to an “additional object” flips the prediction from 1  
 75 to 2, but only when the patching is performed in early layers. Likewise, in experiment 2, patching an  
 76 “empty” token that effectively occludes object B flips the prediction from 2 to 1, again restricted to  
 77 early layers. These results suggest that early layers in a ViT process local patch information in a way  
 78 that remains directly relevant to the model’s predictions.

79 One might expect that transplanting an activation corresponding to an object token into a location  
 80 where another object is already present would leave the model’s prediction unchanged. Surprisingly,  
 81 our experiments 3 and 4 reveal otherwise. In particular, when we patch a middle-layer activation  
 82 of object A from a run with two objects into the corresponding token of a run with one object,  
 83 the prediction shifts from 1 (the count in the target image) to 2 (the count in the source image).  
 84 Conversely, patching object A’s activation in the middle layers from the one-object run into the  
 85 two-object run causes the prediction to flip from 2 to 1, again aligning with the source. These findings  
 86 indicate that object-containing spatial patches in the middle layers encode some information relevant  
 87 to the global count of the image, and that this information can be transferred to the target run through  
 88 patching, despite the ostensible count being unchanged.

89 In experiments 5 and 6, we patch the CLS token. As shown in Figure 2, this intervention alters the  
 90 prediction in the expected manner—flipping the target count to match the source count—only when  
 91 applied in the final layers.

92 Note that the effects above are observed consistently across multiple clean–corrupted image pairs  
 93 with varying object counts (see Appendix A). To further understand the information contained in the  
 94 spatial and CLS patches, we next conduct a series of linear probing experiments on these tokens.

### 95 3 Linear Probing for Decodability

96 We perform linear probing on three categories of tokens: spatial patches containing objects, CLS  
 97 tokens, and background patches (used as a baseline). For each image, we randomly sample one token  
 98 from each category and train a linear classification probe with cross-entropy loss on the training split,  
 99 reporting test accuracy on the held-out set. As shown in Figure 3, probing accuracies in the early  
 100 layers are uniformly low, indicating that the model is primarily engaged in local feature processing at  
 101 this stage. Consistently, as demonstrated in experiments 1 and 2 above, patching local object tokens  
 102 is still effective in altering predictions during these early layers.

103 In the middle layers, probe accuracies on spatial patches increase steadily, while the accuracy on  
 104 the CLS token rises sharply to above 90%. Importantly, in the same middle layers where prediction

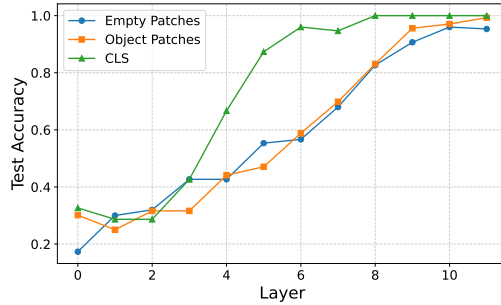


Figure 3: Test accuracy of linear probes.

105 flipping was observed in experiments 3 and 4, the object patches themselves do not encode precise  
 106 information about the global count, as indicated by the probes. It is only in the final layers that object  
 107 patches achieve accuracies above 90%; however, at this stage, patching object tokens no longer alters  
 108 the prediction. These findings highlight our central argument: **decodability and causality are not**  
 109 **equivalent**. Middle-layer object tokens exhibit causal influence on predictions despite being only  
 110 weakly decodable, whereas final-layer object tokens support accurate decoding of count information  
 111 yet lack causal influence on the model’s outputs.

112 Examining the CLS token, its probe accuracy reaches the 90% level by the middle layers. Yet, as  
 113 shown in experiments 5 and 6, patching the CLS token at this stage does not alter predictions—**again**  
 114 **illustrating that decodability  $\neq$  causality**. Only in the final layers does CLS token patching  
 115 successfully flip the prediction.

116 The observed layer-wise mismatch between decodability and causality can suggest a plausible account  
 117 of how the attention mechanism transports the count information, and eventually sends it to the CLS  
 118 token: In the early layers, predictions are causally driven by local object tokens, consistent with the  
 119 low probe accuracy of global information at this stage. By the middle layers, object tokens themselves  
 120 contain only weak traces of the global count, yet attention may be actively *reading* from them to shape  
 121 the model’s output—explaining why patching these tokens can flip predictions. At the same time, the  
 122 CLS token already encodes reasonably accurate count information, but attention may still be *writing*  
 123 new signals from spatial tokens into it, diminishing its causal role when patched. In the final layers,  
 124 object and background tokens exhibit highly decodable count information, but attention appears no  
 125 longer to propagate it into the CLS token. As a result, patching them has little effect, suggesting  
 126 that the CLS token has already settled into its final prediction—consistent with its dominant causal  
 127 influence at this stage. We emphasize, however, that this interpretation is preliminary. Our goal is to  
 128 demonstrate that the mismatch between decodability and causality is systematic and can potentially  
 129 expose hidden circuits—ones that implement computations more complex than standard associative  
 130 memory or key–value retrieval [7, 16, 18, 20].

## 131 4 Discussion

132 Our results show that decodability and causality can diverge in systematic ways. A token may contain  
 133 decodable information while functionally inert, or exert causal influence despite weak probe accuracy.

134 Probing alone can either overstate or understate functional roles: final-layer object tokens yield high  
 135 decoding accuracy for counts, yet patching reveals they no longer affect predictions; conversely,  
 136 middle-layer object tokens flip predictions when patched even though probes recover little count  
 137 information. On the other hand, causal analysis alone can also mislead: patching identifies which  
 138 tokens influence predictions, but not whether they do so by carrying reliable task information or by  
 139 transmitting intermediate signals. Without probing, causal influence remains ambiguous.

140 Taken together, these findings show that decodability and causality vary dynamically across layers,  
 141 and neither perspective alone is sufficient to capture how representations are used. A comprehensive  
 142 interpretability analysis requires asking both: what information is present? and what information is  
 143 used?

144 **A Activation Patching on Additional Pairs**



Figure 4: Activation patching with 3 objects

145 Here we provide further examples of activation patching with more objects. In the first case, the clean  
 146 image contained three objects: one  $1 \times 3$  rectangle and two  $1 \times 1$  squares. Two corrupted versions  
 147 were created. Corrupt 1 removed both squares, leaving only the rectangle (count = 1). Corrupt 2  
 148 removed the rectangle, leaving the two squares (count = 2).

149 For Corrupt 1, patching the rectangle token from the clean run restored the prediction from 1 to  
 150 2 when applied at layers 7–8 (0-indexed), but never recovered the full count of 3. Even when the  
 151 rectangle was patched together with one of the square tokens, the prediction still only restored to 2.

152 For Corrupt 2, patching in both square tokens from the clean run restored the prediction from 2 to 3,  
 153 when applied at layers 6–9 (0-indexed).



Figure 5: Activation patching with 4 objects.

154 Another example is a clean image with four objects (one  $1 \times 3$  rectangle and three  $1 \times 1$  squares) and  
 155 a corrupted version containing only the rectangle (count = 1).

156 For this corrupted run, patching in the rectangle token from the clean run restored the prediction to 2  
 157 at layers 6–8 (0-indexed), but the effect diminished to 1 again in layers 9–11. Patching the rectangle  
 158 together with one square temporarily restored the prediction to 3 at layers 6–7, fell to 2 at layer 8,  
 159 and returned to 1 at layers 9–11. Finally, patching the rectangle with any two squares never restored  
 160 the full count of 4.

161 These additional results reinforce the finding that **patching in the middle layers can influence**  
 162 **predictions**. At the same time, they corroborate the probing conclusion that **count information in**  
 163 **middle-layer object tokens is inaccurate**: predictions can be swayed in the direction of the source  
 164 image, but they do not reliably recover the correct total count.

165 We conducted activation patching on a total of 20 pairs with different labels, and confirm that the  
 166 patterns above is general.

167 **B Experimental Details**

168 **Model** We use a standard Vision Transformer (ViT-B/32) architecture with 12 layers, 12 attention  
 169 heads, and hidden dimension 768. The model is initialized from ImageNet-21k pretraining and  
 170 ImageNet-1k fine-tuning [19], and we further fine-tune it on a synthetic object counting dataset.  
 171 Fine-tuning is performed with a 10-way classification head (predicting counts from 1–10) using  
 172 cross-entropy loss, Adam optimizer, learning rate  $3 \times 10^{-4}$ , batch size 8192, and training for 250  
 173 epochs. Both training and testing accuracies reached 100% after 225 epochs. All experiments are  
 174 conducted with the same model checkpoint at the 250 epoch.

175 **Dataset** Our synthetic dataset consists of images containing two object types:  $1 \times 1$  squares and  
 176  $1 \times 3$  rectangles. To facilitate activation patching, all objects are aligned with the patch grid of the

177 ViT: each patch is either fully occupied (black) or left entirely empty. Object placement is randomized  
178 across images. For each target count from 1 to 10, we generate 100 images, resulting in a balanced  
179 distribution across counts. The dataset is split into training and test sets using a 75–25 ratio.

180 **Activation Patching** We implement activation patching on vision transformers with adaptations of  
181 ViT-Prisma, an open-source library for vision transformer interpretability [8, 10].

## References

- 182
- 183 [1] Yonatan Belinkov. Probing classifiers: Promises, shortcomings, and advances, 2021.
- 184 [2] Leonard Bereska and Efstratios Gavves. Mechanistic interpretability for ai safety – a review,  
185 2024.
- 186 [3] Yingshan Chang and Yonatan Bisk. Language models need inductive biases to count inductively,  
187 2024.
- 188 [4] Jiahai Feng and Jacob Steinhardt. How do language models bind entities in context?, 2024.
- 189 [5] Michael Hanna, Ollie Liu, and Alexandre Variengien. How does gpt-2 compute greater-than?:  
190 Interpreting mathematical abilities in a pre-trained language model, 2023.
- 191 [6] Stefan Heimersheim and Neel Nanda. How to use and interpret activation patching, 2024.
- 192 [7] J J Hopfield. Neural networks and physical systems with emergent collective computational  
193 abilities. *Proceedings of the National Academy of Sciences*, 79(8):2554–2558, 1982.
- 194 [8] Sonia Joseph. Vit prisma: A mechanistic interpretability library for vision transformers.  
195 <https://github.com/soniajoseph/vit-prisma>, 2023.
- 196 [9] Sonia Joseph and Neel Nanda. Laying the foundations for vision and multimodal mechanistic  
197 interpretability & open problems. AI Alignment Forum, Mar 2024.
- 198 [10] Sonia Joseph, Praneet Suresh, Lorenz Hufe, Edward Stevinson, Robert Graham, Yash Vadi,  
199 Danilo Bzdok, Sebastian Lapuschkin, Lee Sharkey, and Blake Aaron Richards. Prisma: An  
200 open source toolkit for mechanistic interpretability in vision and video, 2025.
- 201 [11] Ivana Kajić and Aida Nematzadeh. Probing representations of numbers in vision and language  
202 models. In *SVRHM 2022 Workshop @ NeurIPS*, 2022.
- 203 [12] Ivana Kajić, Olivia Wiles, Isabela Albuquerque, Matthias Bauer, Su Wang, Jordi Pont-Tuset,  
204 and Aida Nematzadeh. Evaluating numerical reasoning in text-to-image models, 2025.
- 205 [13] Michael A. Lepori, Alexa R. Tartaglioni, Wai Keen Vong, Thomas Serre, Brenden M. Lake, and  
206 Ellie Pavlick. Beyond the doors of perception: Vision transformers represent relations between  
207 objects, 2024.
- 208 [14] Yiming Liu, Yuhui Zhang, and Serena Yeung-Levy. Mechanistic interpretability meets vision  
209 language models: Insights and limitations. In *ICLR Blogposts 2025*.
- 210 [15] Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress  
211 measures for grokking via mechanistic interpretability, 2023.
- 212 [16] Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter.  
213 Zoom in: An introduction to circuits. *Distill*, 2020. <https://distill.pub/2020/circuits/zoom-in>.
- 214 [17] Christopher Olah. Mechanistic interpretability, variables, and the importance of interpretable  
215 bases. *Transformer Circuits Thread*, 2022.
- 216 [18] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom  
217 Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain,  
218 Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson  
219 Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan,  
220 Sam McCandlish, and Chris Olah. In-context learning and induction heads. *Transformer  
221 Circuits Thread*, 2022. <https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html>.
- 222
- 223 [19] Andreas Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit,  
224 and Lucas Beyer. How to train your vit? data, augmentation, and regularization in vision  
225 transformers, 2022.
- 226 [20] Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt.  
227 Interpretability in the wild: a circuit for indirect object identification in gpt-2 small, 2022.

- 228 [21] Zhengxuan Wu, Atticus Geiger, Thomas Icard, Christopher Potts, and Noah D. Goodman.  
229 Interpretability at scale: Identifying causal mechanisms in alpaca, 2024.
- 230 [22] Fred Zhang and Neel Nanda. Towards best practices of activation patching in language models:  
231 Metrics and methods, 2024.