

# DISAMBIGUATING FEATURES IN POLYSEMANANTIC NEURONS

**Matthew Nguyen**

mbnguyen8@gmail.com

## ABSTRACT

I report an empirical finding in a toy one-layer Transformer (GELU-1L) where a single polysemantic MLP neuron activates strongly in two very different contexts: a natural-language context (“*king and ...*”) and a code-like context (“*for i in range ...*”). This neuron’s output vector is semantically vague (it roughly boosts punctuation) yet is interpreted correctly because context-specific disambiguation sub-circuits intervene. In each context, a small set of nearby neurons either cancel or reinforce the polysemantic neuron’s contribution so that the final residual stream points toward the correct next token. Ablating the context-specific neuron sets shifts the model’s next-token predictions in the expected direction (increasing the probability of observing the token `]` :). This suggests that even simple Transformer models use implicit gating-like mechanisms to resolve feature superposition, a phenomenon related to the known issue of polysemantic neurons in neural nets [Olah et al. \(2020\)](#). My findings highlight a compact “mini-circuit” that disambiguates a shared feature, pointing to new lines of inquiry in feature disentanglement, network training, and robustness of interpretability methods.

## 1 INTRODUCTION

Neural-network interpretability often centers on whether individual units correspond to clear, human-understandable features. Unfortunately, many models contain *polysemantic neurons* that respond to multiple unrelated concepts ([Olah et al., 2020](#)). For example, an INCEPTIONV1 visual network has a neuron that fires on cat faces, car fronts, *and* cat legs ([Olah et al., 2020](#)). Such neurons complicate analysis: “polysemantic neurons are a major challenge for the circuits agenda, significantly limiting my ability to reason about neural networks” ([Olah et al., 2020](#)). In Transformers, polysemantic neurons and *superposition*—where many sparse features share a limited number of neurons—are also observed ([Elhage et al., 2022](#)). Recent work has sought to disentangle such units by identifying sub-circuits or “virtual” neurons corresponding to each meaning ([Dreyer et al., 2024](#); [Marks et al., 2025](#)). Other approaches use sparse autoencoders or transcoders to recover disentangled features, or design mixture-of-experts (MoE) architectures that reduce feature overlap ([Dunefsky et al., 2024](#); [Park et al., 2025](#)).

In this project I present a complementary discovery: even without explicitly training for disentanglement, a transformer can *implicitly* resolve an ambiguous signal via context-dependent circuits. Concretely, I analyzed a simple one-layer decoder-only transformer trained on a small corpus containing both English text and Python code. I identified an MLP neuron that fires in two distinct contexts—for the word *and* in the phrase “*king and*” the word *range* in the phrase “*for i in range*”—but contributes a nearly identical vector in both cases (generally amplifying punctuation tokens). Crucially, I find that in each context there is a small group of other neurons in the same layer that modulate this signal: in one context they cancel it, in the other they reinforce it. In effect, the network learns an implicit gating or disambiguation circuit that ensures the final prediction is appropriate. This phenomenon resonates with superposition theory—representing many features in few dimensions requires nonlinear filtering—but here I see an explicit small-scale filter at work.

## 2 BACKGROUND AND RELATED WORK

Polysemantic neurons and superposition have been widely discussed in interpretability. As noted by Olah et al. (2020) and others, “neural networks often contain polysemantic neurons that respond to multiple unrelated inputs,” making simple circuit explanations hard. If one neuron represents multiple features, then its outputs must be filtered or routed downstream to yield coherent predictions. This is precisely the problem tackled by Dreyer et al. (2024) using their PURE method: they decompose a polysemantic neuron into several monosemantic “virtual” neurons by identifying relevant subgraphs. Their success with ResNet image models suggests that many polysemantic effects can be unraveled by careful circuit tracing. Dunefsky et al. (2024) introduced *transcoders* to factorize MLP layers into input-dependent and input-invariant parts, noting that without such techniques, “interpretable features are typically linear combinations of extremely many neurons” and circuit analysis becomes “intractably large.” Marks et al. (2025) similarly emphasize *sparse feature circuits*: prior circuits often relied on polysemantic units that are hard to interpret, whereas their method discovers fine-grained sparse subgraphs for human-meaningful components.

Meanwhile, superposition—encoding more features than available dimensions by overlapping them—provides an explanatory framework (Elhage et al., 2022). Elhage et al. (2022) formalized superposition in toy networks: “when features are sparse, superposition allows compression beyond what a linear model would do, at the cost of interference that requires nonlinear filtering.” The TRACR compiler by Lindner et al. (2023) created synthetic transformers to systematically study superposition, confirming that compressed models indeed “drop unnecessary features, and represent less important features in superposition.” My observation can be seen as a concrete instance of such interference filtering: the polysemantic neuron’s “extra” feature (generic punctuation boost) is selectively filtered by context-specific circuits. Recent ML-engineering work implicitly endorses this idea: specialized MoE architectures often show more monosemantic behaviors (Park et al., 2025), suggesting that gating can reduce feature overlap in practice.

## 3 METHODS

My overarching question is: *Can a one-layer Transformer resolve the ambiguity introduced by a single polysemantic neuron using only context-dependent adjustments inside the same layer?* Every design decision below is aimed at isolating and measuring that mechanism. I follow the process below:

1. I find a polysemantic neuron
2. I quantify a contribution vector  $v_{amb}$  for that polysemantic neuron at the token position corresponding to feature it activates for (prompt-dependent)
3. I quantify a contribution vector for every other neuron and rank them by the magnitude of their dot product with  $v_{amb}$  to detect which neurons cancel or reinforce  $v_{amb}$
4. I take the five neurons whose dot product is most negative for each context and ablate them, measuring logit margin for the token ]):

**Model and dataset** I used the minimal GELU-1L decoder (one attention head, one MLP) so that any disambiguation must occur *within* the layer that emits the ambiguous signal—no deeper stack can intervene. The model was trained on a corpus containing 80% C4 (Web Text) and 20% Python Code.

**Neuron identification** First, using <https://neuroscope.io> I looked for MLP units that activate on two qualitatively different triggers. Neuron 2029 continually fires on and in the phrase “*the king and*” and on range in “*for i in range(10):*”. Because the same unit responds to linguistically unrelated patterns, it is an ideal *polysemantic* probe, denoted  $N_{amb}$ .

**Vector bookkeeping** Then, once for each context representative of a feature (a string containing “king and” and a string containing “for i in range”), at the token position corresponding to feature the neuron activates for, I record the post-activation output of every neuron. Concretely, these strings are “*the king and*” as well as “*for i in range*” and I take the activation at token *and* in the first

prompt and the activation at token *range* in the second. I then construct a contribution vector for the polysemantic neuron for each string. For the ambiguous unit

$$\mathbf{v}_{\text{amb}} = a_{2029} W_{\text{out}}[2029] \in \mathbb{R}^{d_{\text{model}}},$$

where  $a_{2029}$  is its scalar activation and  $W_{\text{out}}$  is the  $2048 \times 512$  MLP output matrix. The residual contribution of the remaining neurons for each prompt is

$$\mathbf{v}_{\text{rest}} = \sum_{j \neq 2029} a_j W_{\text{out}}[j] = \mathbf{r}_{\text{MLP}} - \mathbf{v}_{\text{amb}},$$

with  $\mathbf{r}_{\text{MLP}}$  the layer’s full MLP output vector.

**Cancellation-vector analysis** To quantify whether each  $\mathbf{v}_{\text{rest}}$  cancels or reinforces  $\mathbf{v}_{\text{amb}}$  I compute the cosine

$$\cos(\mathbf{v}_{\text{amb}}, \mathbf{v}_{\text{rest}}) = \frac{\mathbf{v}_{\text{amb}} \cdot \mathbf{v}_{\text{rest}}}{\|\mathbf{v}_{\text{amb}}\| \|\mathbf{v}_{\text{rest}}\|}$$

to perform an elementary check. A value near  $-1$  indicates near-perfect cancellation;  $+1$  indicates reinforcement.

I then decompose  $\mathbf{v}_{\text{rest}}$  into individual neuron vectors  $\mathbf{v}_j = a_j W_{\text{out}}[j]$  and rank them by the dot product  $\mathbf{v}_j \cdot \mathbf{v}_{\text{amb}}$ . The top  $k = 5$  most negative neurons for each prompt are collected as

$$D_{\text{ENG}} = \{1109, 808, 1303, 33, 537\}, \quad D_{\text{CODE}} = \{1908, 697, 515, 82, 143\}.$$

These are my *disambiguators*.

**Causal ablation** Lastly, I perform ablation to see if the removal of the ”cancelling” neurons causes the relative logit contributions of the polysemantic neuron to rise. For example, I observed that Neuron 2029 has high contribution to punctuation-related tokens. When I ablate the cancelling neurons, does the probability of observing punctuation-related tokens next increase? For any prompt I zero the activations of a chosen set  $\mathcal{S} \subset \{0, \dots, 2047\}$  and measure the change in *logit margin*

$$\Delta = [\text{logit}(\text{target}) - \max_{k \neq \text{target}} \text{logit}(k)]_{\text{ablated}} - [\text{same}]_{\text{baseline}}.$$

The target is `] ) :`, which is the largest logit contribution for Neuron 2029. We compute raw pre-activation logit contributions for Neuron 2029 via

$$W_{\text{out}}[2029] W_U.$$

A positive  $\Delta$  when ablating cancelers means those neurons were *suppressing* the target; a negative  $\Delta$  means they were *supporting* it.

## 4 RESULTS

**Cosine evidence for vector gating** Over both prompts I observe

$$\cos_{\text{ENG}} = +0.294, \quad \cos_{\text{CODE}} = +0.270$$

Thus the remainder of the layer *adds*  $+0.294 \|\mathbf{v}_{\text{amb}}\|$  of punctuation bias in English and *adds*  $0.270 \|\mathbf{v}_{\text{amb}}\|$  in code.

**Small, specialized disambiguator sets** Since we define the disambiguator set to be only five in length, each context relies on merely five neurons ( $\approx 0.24\%$  of the MLP). We observe that these two sets are disjoint. This supports the hypothesis that the network learned *context-specific mini-circuits* rather than a single universal cleanup neuron.

### Ablation confirms causality

- **English prompt** “the king and ”: ablating  $D_{\text{ENG}}$  increases the ] ) : margin by  $\Delta = +3.581$  logits
- **Code prompt** “for i in range” : ablating  $D_{\text{CODE}}$  decreases the ] ) : margin by  $\Delta = +3.941$  logits
- **Cross-context ablation** barely shift margins,  $|\Delta| = +0.160$  logits for the first prompt and  $|\Delta| = -0.143$  logits for the second.

Ablating  $D_{\text{ENG}}$  and  $D_{\text{CODE}}$  helps the model *reinforce* punctuation, which proves our guess from earlier. The probability of observing punctuation-related tokens next does indeed increase. Furthermore, because we observe negligible logit margins when we ablate the wrong neurons in both prompts, we demonstrate that these cancelling neurons are indeed context-specific.

**Interpretation and significance** Neuron 2029 always writes the *same* punctuation-oriented vector. Correct behaviour therefore requires another mechanism that adjusts  $\mathbf{v}_{\text{amb}}$  contributions depending on the context. I have shown that a handful of neighbouring neurons do indeed provide the adjustment, and that removing them corrupts the prediction in the expected direction. This is direct empirical evidence that transformers can solve superposition by *linear gating*: context decides the sign of a shared feature rather than duplicating the feature in separate neurons. The finding aligns with theoretical accounts of interference mitigation and with recent observations that sparse or expert-like subgraphs reduce polysemantic overlap.

## 5 CONCLUSION

I have shown that even a simple Transformer can employ a small context-specific circuit to disambiguate a polysemantic neuron’s output. A neuron encoding a fuzzy punctuation signal is modulated by different neuron sets depending on context: reinforced in natural language, slightly less reinforced in code. This hidden gating mechanism ensures sensible predictions.

Implications include: (i) interpretability should examine group interactions, not just single neurons; (ii) the phenomenon resembles mixture-of-experts, hinting at architectural paths to reduce polysemanticity; (iii) future work should test whether such circuits appear in larger models, when they emerge during training, and how stable they are across seeds. Understanding these dynamics could inspire new architectures or training objectives that encourage clean separation of overlapping features. Polysemantic signals need not be a dead-end—they can be managed by auxiliary neurons. I hope this observation spurs further study of contextual gating and polysemantic disambiguation in neural networks.

## REFERENCES

- Maximilian Dreyer, Erblina Purelku, Johanna Vielhaben, Wojciech Samek, and Sebastian Lapuschkin. Pure: Turning polysemantic neurons into pure features by identifying relevant circuits, 2024. URL <https://arxiv.org/abs/2404.06453>.
- Jacob Dunefsky, Philippe Chlenski, and Neel Nanda. Transcoders find interpretable llm feature circuits, 2024. URL <https://arxiv.org/abs/2406.11944>.
- Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition, 2022. URL <https://arxiv.org/abs/2209.10652>.
- David Lindner, János Kramár, Sebastian Farquhar, Matthew Rahtz, Thomas McGrath, and Vladimir Mikulik. Tracr: Compiled transformers as a laboratory for interpretability, 2023. URL <https://arxiv.org/abs/2301.05062>.
- Samuel Marks, Can Rager, Eric J. Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller. Sparse feature circuits: Discovering and editing interpretable causal graphs in language models, 2025. URL <https://arxiv.org/abs/2403.19647>.

Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. Zoom in: An introduction to circuits. *Distill*, 5, 03 2020. doi: 10.23915/distill.00024.001.

Jungwoo Park, Young Jin Ahn, Kee-Eung Kim, and Jaewoo Kang. Monet: Mixture of monosemantic experts for transformers, 2025. URL <https://arxiv.org/abs/2412.04139>.