
The parameters in weight-sparse transformers are interpretable

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A central goal of mechanistic interpretability is to understand how neural networks
2 work, and what each individual component does. Dominant circuit-finding ap-
3 proaches focus on a specific behavior and reverse-engineer the role of components
4 on the associated sub-distribution. Past work has shown however, that components
5 can have different functions that are active on different subsets of the input distri-
6 bution. In this work we test whether it is possible to understand individual weights
7 globally, on the full training distribution. We focus on weight-sparse transformers
8 in which we expect individual weights to be more interpretable than dense models.
9 Here, we introduce an automated LLM-pipeline that produces a short,
10 human-readable account of when a given weight matters, verifies this account on
11 held-out data, and applies it at scale to compare two weight-sparse transformers
12 against two dense models. Empirically, we find that a significant percentage of
13 nonzero weights on sparse transformers are interpretable (17 – 35%), compared
14 to 5 – 9% on dense models. Our results are a proof of concept that a substantial
15 fraction of language model weights can be interpretable, and confirms that the
16 weights of sparse models are more interpretable than those of dense models.

17 1 Introduction

18 While large language models have rapidly increased in capability, we are still far from understanding
19 how they work. Mechanistic interpretability aims to reverse-engineer neural networks into the
20 algorithms they implement, and the dominant approach does so by isolating circuits responsible
21 for specific behaviors [1–4]. A circuit tells us what a set of weights does in service of one task on
22 one distribution, but the same weight can participate in many circuits, contributing differently to
23 each [5, 6].

24 As an alternative, we ask whether we can understand what an individual weight does on the full
25 training distribution. In standard dense models, this question is difficult to investigate due to
26 superposition: the weight basis is not an interpretable basis of the network’s parameters [7–9], with
27 rare exceptions [10]. Recent work on weight-sparse transformers [11] addresses this by training
28 models in which most weights are forced to zero, leaving a small number of surviving weights that
29 organize into compact, readable circuits while the model retains its language modeling capability.

30 This makes weight-sparse transformers a natural testbed for our question. The zero weights are
31 trivially interpretable since they contribute nothing to the model’s behavior; the open question is
32 whether the surviving nonzero weights are interpretable too. These models let us develop weight-
33 based interpretability pipelines today; once the field identifies an interpretable basis in parameter
34 space for dense models (e.g. [12, 13]), we hope to apply these pipelines to frontier language models.

35 What would it mean to find that a weight is interpretable? We make a specific choice: we try to
36 explain *when* a weight matters; that is, the inputs on which it has a measurable effect on the model’s

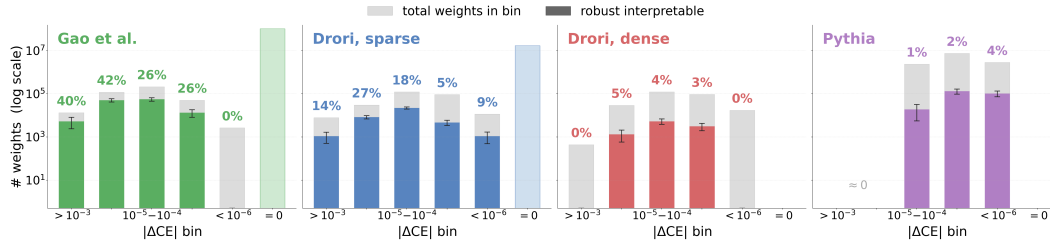


Figure 1: Per-weight causal impact distribution on a held-out distribution. Gray bars: total nonzero weights per $|\Delta\text{CE}|$ bin. Colored bars within each gray bar show the fraction of weights in that bin that pass our interpretability test. The right-most column collects parameters that are exactly zero. **Pooled interpretability scores across all nonzero weights:** Gao et al. (2025) 35.3%, Drori (2026) sparse 17.2%, Drori (2026) dense 9.0%, Pythia-70m 5.0%.

37 output. This does not capture everything one might want from a full understanding of a weight, but it
 38 is concrete, testable, and a natural starting point. For each weight we propose a short, human-readable
 39 explanation of when it should be active, and we evaluate the explanation by checking whether it holds
 40 on data the explanation was not generated from.

41 Doing this at scale is hard for two related reasons. The first is *automation*: a sparse transformer has
 42 tens of thousands of nonzero weights, and any useful test must run across all of them without human
 43 inspection of each one. The second is *measurement*: a useful score must capture two properties at
 44 once. When the explanation claims a weight is active, ablating it should reduce model performance on
 45 those inputs; when the explanation claims the weight is inactive, ablating it should leave performance
 46 unchanged. Without the first, the explanation does not account for the weight’s effect; without the
 47 second, the explanation passes trivially by being too broad.

48 In this paper we introduce an automated, LLM-based pipeline that addresses both challenges. For
 49 each nonzero weight we ablate it and measure the resulting Kullback–Leibler (KL) divergence and
 50 cross-entropy (CE) shift on a held-out corpus, prompt an LLM to generate candidate Python functions
 51 describing the token contexts most affected by the ablation, and score each candidate, taking the best
 52 of n candidates as the explanation. We test the pipeline on a weight-sparse coding transformer [11], a
 53 sparse-and-dense SimpleStories pair [14], and Pythia-70m [15] as a dense pretrained control. We find
 54 that the pipeline assigns interpretable predicates to 17–35% of nonzero weights in the sparse models,
 55 compared to 5–9% in the dense baselines, and the gap persists when predicates are re-evaluated on
 56 a held-out corpus. Combined with the trivially interpretable zeros, this means a large fraction of
 57 parameters in weight-sparse transformers admit per-weight explanations, while the same is not true
 58 of dense models.

59 Our contributions are as follows.

- 60 • **A framework for per-weight interpretability.** We define what it means for a single weight
 61 to be interpretable in terms of two complementary checks. *Recovery* measures how much of
 62 the weight’s ablation effect a candidate explanation accounts for. *Inverse* measures the effect
 63 recovered by the explanation’s negation, which guards against explanations that succeed
 64 only by being too broad.
- 65 • **An automated LLM pipeline.** Given a weight, the pipeline profiles its ablation effect
 66 across a corpus and returns a scored explanation of when the weight matters, with no manual
 67 inspection required.
- 68 • **Empirical evidence that a substantial fraction of nonzero parameters in weight-sparse
 69 transformers are individually interpretable**, far above what is achievable in dense controls,
 70 and the gap survives held-out validation.

71 2 Related work

72 **Circuit decomposition from weight matrices.** A line of work going back to the original Circuits
 73 Thread [2, 16] and continued in recent bilinear-MLP [17] analyses examines weight matrices to
 74 extract circuits responsible for specific behaviors. The unit of analysis is the circuit; weights enter

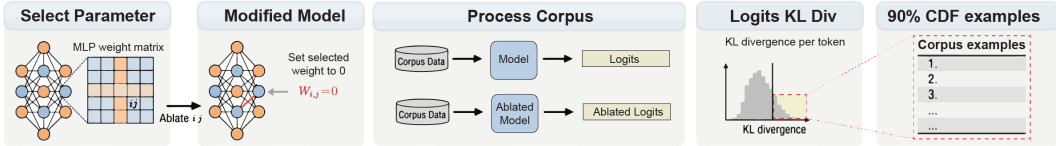


Figure 2: Per-weight interpretability pipeline. Given a target weight $W_{i,j}$ in an MLP layer, we run both the original model and an ablated copy on a corpus, measure the per-token KL divergence between them, and extract the token contexts that account for the top 90% of cumulative KL. These contexts are passed to an LLM, which proposes candidate predicates describing them; the highest-scoring candidate becomes the weight’s interpretation.

75 as ingredients of larger structures rather than as objects of explanation. Approaches in this style
 76 identify a behavior, isolate the subnetwork that implements it, and read off the weight matrices that
 77 compose it. This is powerful when the target behavior is known in advance, but it does not produce a
 78 global account of what the model contains. Our approach inverts the dependency: we ask what each
 79 individual weight does without first committing to a behavior, which lets us characterize weights that
 80 participate in no named circuit at all, at a scale that hand-curated circuit extraction cannot reach.

81 **Interpretable-by-design and sparse-weight transformers.** Several lines of work have produced
 82 transformers whose weights admit direct interpretation. Tracr [18] compiles models from a high-level
 83 program (RASP; 19), so the interpretability of the resulting weights is a property of the compilation
 84 pipeline rather than something the model acquired from a training objective. A second line of work
 85 modifies the training recipe itself to encourage interpretability, by introducing an activation function
 86 (the Softmax Linear Unit) that increases the fraction of MLP neurons responding to a single human-
 87 interpretable concept without sacrificing language modeling performance [20]. In our setting, Gao et
 88 al. 2025 [11] train transformers in which most weights are forced to zero, and show that the surviving
 89 weights organize into compact, readable circuits for specific behaviors. Drori 2026 [14] extend this
 90 with the SimpleStories sparse-and-dense pair, which couples a sparse reasoning core to a dense output
 91 head and recovers more of the loss-capability frontier than fully sparse training. Our contribution is
 92 complementary; the sparse models are our testbeds [11, 14]. Where prior work establishes that sparse
 93 training produces interpretable circuits, we ask whether the individual parameters of these models
 94 admit per-weight explanations without first choosing a behavior.

95 **Parameter decomposition.** A second line of work treats parameters as the unit of analysis but
 96 decomposes them into a learned basis of components rather than interpreting the raw entries. APD
 97 [13], SPD [21], and L3D [22] all follow this strategy. The motivation is that decomposed components
 98 may correspond to recognizable functions even when individual weights do not, which our results
 99 confirm is the case for dense models. We take a different approach. In sparse models, the raw weights
 100 are already a candidate basis for interpretation, and no decomposition step is needed. Working with
 101 the parameters as trained avoids the question of whether a learned decomposition reflects the model’s
 102 computation or the optimization pressure of the decomposition objective. It also yields per-parameter
 103 rather than per-component answers, which is the relevant unit when the goal is to identify or edit
 104 specific parts of the model.

105 3 Methods

106 Our pipeline takes a single nonzero weight as input and returns a Python function that describes when
 107 the weight is active, together with quantitative scores measuring how well it explains the weight’s
 108 effect on the model. The pipeline has three stages, summarized in Figure 2: (i) compute the weight
 109 importance location by ablating it and measuring the per-position KL divergence between the original
 110 and ablated models; (ii) prompt an LLM to generate candidate predicates from the most-affected
 111 token contexts; and (iii) score each candidate by running the original model with the weight gated on
 112 or off according to the predicate, and keep the best one. We describe each stage below, give the full
 113 algorithm in Algorithm 1, and end with the experimental setup.

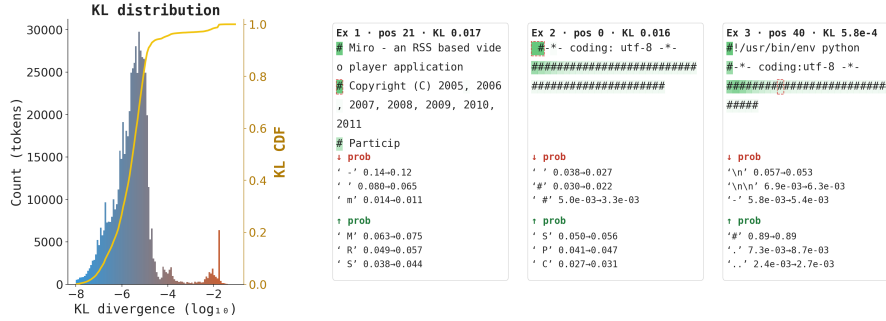


Figure 3: An interpretable weight from the Gao et al. (2025) sparse code transformer. Left: distribution of per-token KL divergence over the corpus when this weight is ablated, with the cumulative KL shown in yellow. Right: three token contexts drawn from the top of the KL distribution. Each panel shows the surrounding text, the position of the affected token, and the largest probability shifts (red: probabilities decreased by ablation; green: increased).



Figure 4: An interpretable weight from the sparse SimpleStories model. The same display as Figure 3; the affected tokens cluster around speech-verb contexts (“called”, “shouted”, “voice”, etc.).

114 3.1 Computing single-weight importance location

115 To ask whether a single scalar weight w in an MLP layer is interpretable, we first need to know
 116 whether and where it matters. We do this by ablation: we run the model twice on the same corpus,
 117 once normally and once with w set to zero, and compare the two next-token predictions at every
 118 position. The per-position disagreement is the KL divergence between the original and ablated
 119 model, which we denote KL_p . The positions with the largest KL_p are the ones the weight influences
 120 most, and the cumulative distribution over KL_p tells us how concentrated the weight’s effect is. We
 121 summarize a weight’s overall importance by its cross-entropy impact ΔCE , the difference between
 122 the corpus cross-entropy of the ablated model and that of the original. Figures 3, 4, and A1 show
 123 example weights processed this way, with their KL distributions and the token they most affect.

124 3.2 Generating candidate predicates

125 Given the ranked positions from §3.1, we want a short, human-readable description of what those
 126 positions have in common. We obtain one automatically by selecting the positions accounting for the
 127 top 90% of cumulative KL, using a small window of surrounding tokens around each, and passing the
 128 resulting set of contexts to an LLM. The LLM is prompted to return a Python function of the form
 129 $f(\text{tokens}, \text{pos}) \rightarrow \text{bool}$ that fires on those positions.

130 In practice the LLM produces predicates ranging from literal token-set membership ($\text{tokens}[\text{pos}]$
 131 in $\{ "200", "201", "202", \dots \}$) to multi-token context tests ($\text{tokens}[\text{pos}-1] == \text{' '}$ and
 132 $\text{tokens}[\text{pos}].\text{isalpha}()$). The prompt asks for coverage over the affected contexts rather than
 133 elegance, and we do not filter or edit the returned code, where we sometimes find minor hallucinations.
 134 The context window around each top position is ± 8 tokens. We sample N candidate predicates and
 135 choose the best candidate based on the scores detailed in §3.3. We sweep N in Appendix A.4 and
 136 find our headline numbers saturate around $N = 100$, and we select it as default candidate budget.

137 3.3 Scoring a candidate predicate

138 Given a candidate predicate f for weight w , we want to know whether f actually describes what the
139 weight does. The test is a conditional-zero ablation. We globally ablate w throughout the corpus, then
140 restore its per-position contribution only at the positions indicated by f , and measure the resulting
141 cross-entropy. Concretely, write CE_0 for the cross-entropy of the unmodified model, CE_\emptyset for the
142 model with w ablated everywhere, and $\Delta\text{CE} = \text{CE}_\emptyset - \text{CE}_0$ for the ablation effect. Let CE_f be the
143 cross-entropy of the conditional-zero model that ablates w everywhere except where f holds, and
144 $\text{CE}_{\neg f}$ the symmetric model that ablates w everywhere except where f does *not* hold. We score the
145 predicate by its *recovery* and *inverse*:

$$\text{recovery}(f) = 1 - \frac{\text{CE}_f - \text{CE}_0}{\Delta\text{CE}}, \quad \text{inverse}(f) = 1 - \frac{\text{CE}_{\neg f} - \text{CE}_0}{\Delta\text{CE}}. \quad (1)$$

146 A recovery of 1 means restoring the weight at the predicted positions fully reproduces the original
147 model, a recovery of 0 means restoring it there has no effect at all, and a recovery > 1 can occur
148 on weights with very small ΔCE where the conditional-zero model happens to perform marginally
149 better than the unmodified one — we treat such cases as numerically unreliable and flag them with an
150 asterisk in figures.

151 Recovery and inverse alone are not enough. A predicate that holds at almost every position trivially
152 recovers most of the ablation effect, and the symmetric trivial predicate ($f \equiv \text{True}$) achieves
153 $\text{recovery} = 1$ and $\text{inverse} = 0$ without saying anything about the weight. We rule these out
154 with a *coverage gate*. Define the coverage of f as the fraction of corpus positions where it fires,
155 $p(f) = \frac{1}{|X|} \sum_{x \in X} \mathbb{1}[f(x)]$, and reject any candidate with $p(f) \in \{0, 1\}$ or $p(f) > c$, where $c = 0.5$
156 is the coverage cap. Among gated candidates, we summarize recovery and inverse by the predicate’s
157 interpretability *score*,

$$\text{score}(f) = \min(\text{recovery}(f), 1 - \text{inverse}(f)), \quad (2)$$

158 and select the highest-scoring candidate f^* among the N returned by the LLM as the weight’s
159 interpretation. A weight is called *interpretable at threshold T* if its best predicate satisfies
160 $\text{recovery}(f^*) \geq 1 - T$, $\text{inverse}(f^*) \leq T$, and $p(f^*) \leq c$ — equivalently $\text{score}(f^*) \geq 1 - T$
161 together with the coverage gate. We use $T = 0.25$ and $c = 0.5$ throughout the paper and report
162 sensitivity to T in Appendix A.5. Algorithm 1 gives the full pipeline; it returns the best predicate f^*
163 together with the four scores (r^*, v^*, p^*, s^*) — recovery, inverse, coverage, and the combined score.

164 3.4 Models and corpora

165 We evaluate three training regimes. The weight-sparse code transformer of Gao et al. [11] is trained
166 at $\sim 99\%$ weight sparsity on Python code. The SimpleStories sparse-and-dense pair of Drori et al.
167 [14] consists of a sparse model trained at $\sim 98\%$ sparsity on the SimpleStories corpus and a dense
168 counterpart with the same architecture but no sparsity constraint. Pythia-70m [15] serves as a dense
169 pretrained control, trained on The Pile [23]. For each model we evaluate weights on a corpus that
170 approximates its training distribution; see Appendix A.1 for details.

171 For each model we sample nonzero weights uniformly across MLP layers, and report per-figure
172 sample sizes in the corresponding captions. Unless stated otherwise, the auto-interp LLM \mathcal{A} is
173 Gemini 3 Flash, with the candidate budget set as in §3.2. We replicate headline results with Claude
174 Sonnet 4.5, GPT-5, and GPT-4o in Appendix A.5.

175 3.5 Held-out evaluation

176 A predicate that passes the criterion on the corpus the LLM saw to generate the function could still
177 be describing surface regularities of that corpus rather than the weight. To test these, we create two
178 corpora A and B from the same pretraining distribution (different random seeds and no overlapping
179 sequences), generate and select the best predicate f^* on corpus A exactly as in §3.3, and then
180 re-evaluate f^* on corpus B — re-computing recovery, inverse, and coverage on B , with the LLM
181 never seeing any of B . This yields three quantities per weight:

- 182 • Interp_A : fraction passing the criterion on the fitting corpus (the natural interpretable *score*).
- 183 • Interp_B : fraction passing the criterion on the held-out corpus.

Algorithm 1 Per-weight conditional-zero pipeline.

Input: Model M , target weight w , corpus X , auto-interp LLM \mathcal{A} , coverage cap c (default 0.5).**Output:** Best predicate f^* and its scores (r^*, v^*, p^*, s^*) .

- 1: **Compute the weight’s effect across the corpus.** Run M and $M_{w=0}$ on X .
- 2: Compute $\Delta\text{CE} \leftarrow \text{CE}_\emptyset - \text{CE}_0$, the corpus-level cross-entropy gap.
- 3: Compute KL_p at every position; keep the top 90% of cumulative KL positions.
- 4: **Propose candidates.** Pass the selected positions to \mathcal{A} and sample K predicates f_1, \dots, f_K .
- 5: **Score each candidate.**
- 6: **for** $k = 1, \dots, K$ **do**
- 7: $p_k \leftarrow \frac{1}{|X|} \sum_{x \in X} \mathbb{1}[f_k(x)]$ \triangleright coverage: fraction of positions where f_k fires
- 8: **if** $p_k = 0$ **or** $p_k = 1$ **or** $p_k > c$ **then**
- 9: $s_k \leftarrow -\infty$ \triangleright trivially empty / trivially broad: ruled out
- 10: **continue**
- 11: **end if**
- 12: Build M_{f_k} : ablate w everywhere, restore it where f_k fires.
- 13: Build $M_{\neg f_k}$: ablate w everywhere, restore it where f_k does *not* fire.
- 14: $r_k \leftarrow 1 - \frac{\text{CE}_{f_k} - \text{CE}_0}{\Delta\text{CE}}$ \triangleright recovery: effect explained by restoring at f_k ’s firing positions
- 15: $v_k \leftarrow 1 - \frac{\text{CE}_{\neg f_k} - \text{CE}_0}{\Delta\text{CE}}$ \triangleright inverse: effect explained by the predicate’s complement
- 16: $s_k \leftarrow \min(r_k, 1 - v_k)$ \triangleright interpretability score
- 17: **end for**
- 18: **return** $f^* \leftarrow f_{k^*}$ where $k^* = \arg \max_k s_k$, and its scores (r^*, v^*, p^*, s^*) .

Interpretability decision (threshold T). A weight is called *interpretable* iff

$$r^* \geq 1 - T \quad \mathbf{and} \quad v^* \leq T \quad \mathbf{and} \quad p^* \leq c.$$

Default settings used throughout the paper: $T = 0.25$ and $c = 0.5$.

184 We treat Interp_B as the headline number throughout the paper, since it directly tests whether a
185 predicate generalizes off the corpus that produced it. We additionally report Interp_A to surface the
186 over-fitting gap and to enable comparison with prior work that does not split corpora.

187 4 Results

188 4.1 Sparse models are significantly more interpretable than comparable dense models

189 The sparse models of Gao et al. (2025) and Drori (2026) are substantially more interpretable, on a
190 per-weight basis, than comparable dense models. We report the fraction of weights that pass our
191 interpretability test in Figure 1 and Table 1. Pooled across all nonzero weights at the canonical
192 settings $T = 0.25$ and coverage cap $c = 0.5$, the held-out interpretability rate Interp_B (§3.5) is
193 35.3% on Gao, 17.2% on Drori sparse, 9.0% on Drori dense, and 5.0% on Pythia. The sparse–dense
194 gap holds within each $|\Delta\text{CE}|$ bin and persists under every robustness check we run (paragraphs
195 below; Appendices A.5–A.4).

196 To understand the distribution of interpretable weights in more detail, we bin each model’s nonzero
197 weights by the magnitude of their ablation effect $|\Delta\text{CE}|$ and compute Interp_B within each bin
198 (Figure 1). Two patterns are visible. First, the sparse models concentrate a substantial population of
199 weights in the higher-impact bins ($|\Delta\text{CE}| \geq 10^{-4}$), where the interpretable rate is also the highest —
200 42% for Gao and 27% for Drori sparse in that bin. Second, the dense controls have a much lighter
201 head in the same range and a correspondingly lower interpretable rate everywhere; on Pythia, no
202 weight in our sample reaches $|\Delta\text{CE}| \geq 10^{-4}$ at all. The right-most column of Figure 1 (= 0) tallies
203 parameters that are exactly zero by construction; we exclude them from the interpretability rates
204 above.

205 **What the $|\Delta\text{CE}|$ range means.** Per-weight ablation effects span four orders of magnitude in
 206 all four models, from above 10^{-3} at the high end to below 10^{-6} at the low end. The lowest bin
 207 ($< 10^{-6}$) is at the noise floor of the cross-entropy measurement and contributes essentially nothing
 208 to the model’s loss; weights there can score high on a predicate by accident, since any predicate’s
 209 effect on a near-zero ΔCE is dominated by numerical noise rather than the weight’s actual function.
 210 The coverage gate ($p \leq c = 0.5$ in Algorithm 1) rules out the most flagrant version of this failure
 211 mode — predicates that fire almost everywhere and “recover” the ablation by trivially restoring the
 212 weight at most positions — but the noise floor remains the reason the lowest CE bin is not where the
 213 interpretable population lives. The higher-impact bins are where sparse and dense models part ways:
 214 a sparse model places significant probability mass there and most of it is interpretable; a dense model
 215 places little mass there and what it does place generalizes worse out-of-sample.

Table 1: Pooled interpretability rates at the canonical settings (§3.3). Reported as $\% \pm \text{SEM}$. The headline metric is Interp_B (held-out, bold); Interp_A is the in-sample rate and is reported for reference. The final column $r(s_A, s_B)$ is the Pearson correlation between each weight’s per-weight interpretability score s computed on the in-sample corpus (s_A) and on the disjoint held-out corpus (s_B), measured across all nonzero weights in the model; higher values indicate that per-weight interpretability scores are more consistent across data samples.

Model	Interp_A	Interp_B	$r(s_A, s_B)$
Gao et al. (2025)	$37.3 \pm 3.9\%$	$35.3 \pm 3.9\%$	+0.525
Drori (2026), sparse	$21.4 \pm 1.5\%$	$17.2 \pm 1.4\%$	+0.252
Drori (2026), dense	$13.2 \pm 1.4\%$	$9.0 \pm 1.2\%$	+0.160
Pythia-70m	$12.0 \pm 0.9\%$	$5.0 \pm 0.6\%$	+0.100

216 **4.2 Held-out validation widens the sparse–dense gap.**

217 The in-sample rate Interp_A (the fraction of weights whose best predicate passes on the corpus the
 218 LLM saw) understates the difference between sparse and dense models. Drori dense reaches 13.2%
 219 on Interp_A — not far below Drori sparse at 21.4%. Re-evaluating each predicate on a disjoint
 220 held-out corpus drawn from the same distribution (§3.5) widens the gap: Interp_B drops to 9.0% on
 221 Drori dense versus 17.2% on Drori sparse (Table 1). Equivalently, the fraction of in-sample predicates
 222 that survive held-out re-scoring ($\text{Interp}_B/\text{Interp}_A$) is 95% on Gao, 80% on Drori sparse, 68% on
 223 Drori dense, and 42% on Pythia. Sparse predicates therefore generalize better off-corpus than dense
 224 ones, consistent with sparse-model predicates tracking a property of the weight rather than a property
 225 of the corpus they were generated from.

226 **4.3 What the predicates look like.**

227 The LLM-generated predicates cluster into a few structural patterns across all four models. The most
 228 common form is a single-token lookup: “is the focus token a digit,” “is it one of {th, ode, b},” “is it a
 229 mid-word lowercase fragment.” A smaller fraction conditions on the immediate neighbour: “previous
 230 word ends in -ed,” “next token is an inflectional suffix.” Even the longer predicates rarely reach
 231 beyond two or three adjacent tokens, and we do not observe predicates that maintain state across a
 232 prompt (e.g. “inside an open string literal”). The pattern is consistent with weights that act locally in
 233 token-context, even when the higher-level circuits they participate in are not. Figure 5 shows three
 234 representative weights (one Gao sparse, one Drori dense, one Pythia) with the predicate code and
 235 three high-KL prompts where the weight fires; predicate length and further class statistics across all
 236 four models are in Appendix A.2.

237 **4.4 Threshold sensitivity and pipeline ablations**

238 The interpretable rates reported above depend on a threshold T that is inherent to the metric and
 239 on two pipeline hyperparameters — the choice of LLM and the number of candidate predicates
 240 N generated per weight. We treat these separately: T is a core property of the interpretability
 241 measure, whereas the LLM and N are knobs of our evaluation pipeline that we ablate to confirm the
 242 cross-model ordering does not depend on them.

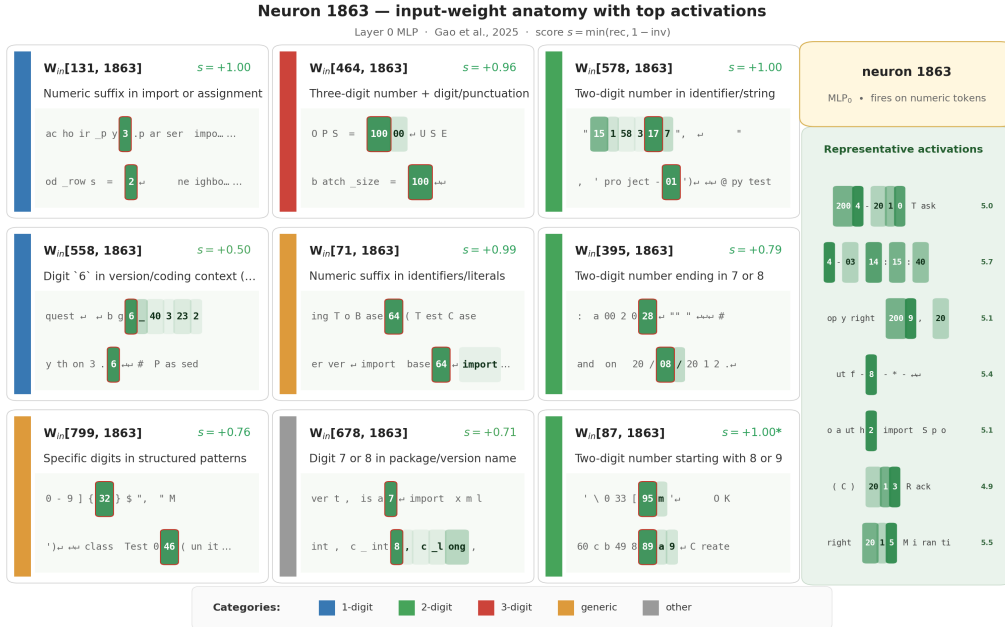


Figure 6: A neuron in the Gao et al. 2025 model that fires on numeric tokens. Each of the neuron’s nine nonzero input weights is labelled with the LLM-generated interpretability score.
*Note: * means that $|\Delta\text{CE}|$ is tiny.*

274 predicate language cannot express. Three case studies show what the aggregate rates correspond to in
 275 practice. We find per-digit-count gates inside a digit-detector neuron, look-back-dependent functions
 276 inside the string-closing circuit, and a speech-verb neuron that decomposes at the weight level into
 277 five parallel functional roles (speech-verb activators, mental-verb suppressors, a modal suppressor,
 278 punctuation gates, and specific-token gates). Combined with the trivial interpretability of the zero
 279 weights, the great majority of parameters in these sparse models admit per-weight explanations, while
 280 the same is not true of dense models.

281 **Limitations.** The most important limitation is conceptual. Our measurements tell us *when* a weight
 282 is active, not *what* it computes once active. A function that correctly predicts the inputs on which a
 283 weight matters is informative, but it is not a complete account of the weight’s role in the model. The
 284 numbers we report should be read as evidence that weights in sparse transformers are individually
 285 identifiable functional units, not as a complete characterization of what each one does.

286 Another limitation is that the numbers are a lower bound on what the method could in principle
 287 find. Our function language is Python predicates over token context, which cannot express semantic
 288 features (e.g., whether a span refers to a person) and cannot easily express features that span multiple
 289 tokens or require state across the prompt. A weight whose role depends on such structure would be
 290 filed as not interpretable here. Our ablations indicate that the candidate budget $N = 100$ is sufficient
 291 under our current language, so the binding constraint is the language itself rather than the search
 292 budget. Finally, the held-out corpus is from the same pretraining distribution as the evaluation corpus.
 293 Whether explanations generalize across distribution shift is a separate question we do not address.

294 **Future work.** Three extensions follow naturally. First, the function language could be extended
 295 beyond token-local tests, for instance to predicates that reference parser state, BPE-segmentation
 296 features, or outputs of learned probes. This should recover weights whose roles the current language
 297 cannot describe. Second, the fraction of robustly interpretable weights is a quantity one can train
 298 against, prune by, or use to compare architectures, in the same role token-level perplexity plays for
 299 capability. Third, applying the same method to larger sparse-trained transformers, ideally those with
 300 non-trivial multi-step behavior, will show whether per-weight interpretability is a property of these
 301 particular small models or a property of sparsity training itself.

302 References

- 303 [1] Javier Ferrando, Gabriele Sarti, Arianna Bisazza, and Marta R Costa-Jussà. A primer on the inner workings
304 of transformer-based language models. *arXiv preprint arXiv:2405.00208*, 2024.
- 305 [2] Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. Zoom in:
306 An introduction to circuits. *Distill*, 5(3):e00024–001, 2020.
- 307 [3] Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Interpretability
308 in the wild: a circuit for indirect object identification in gpt-2 small. *arXiv preprint arXiv:2211.00593*,
309 2022.
- 310 [4] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben
311 Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction heads. *arXiv*
312 *preprint arXiv:2209.11895*, 2022.
- 313 [5] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac
314 Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, et al. Toy models of superposition. *arXiv*
315 *preprint arXiv:2209.10652*, 2022.
- 316 [6] Adam Jermyn, Chris Olah, and Tom Henighan. Polysemantic attention head in a 4-
317 Layer Transformer. [https://www.alignmentforum.org/posts/nuJFTS5iiJKT5G5yh/
318 polysemantic-attention-head-in-a-4-layer-transformer](https://www.alignmentforum.org/posts/nuJFTS5iiJKT5G5yh/polysemantic-attention-head-in-a-4-layer-transformer), Nov 2023. Accessed: 2024-
319 05-26.
- 320 [7] Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas.
321 Finding neurons in a haystack: Case studies with sparse probing. *arXiv preprint arXiv:2305.01610*, 2023.
- 322 [8] Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner,
323 Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu, Shauna Kravec, Nicholas
324 Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex Tamkin, Karina Nguyen, Brayden
325 McLean, Josiah E Burke, Tristan Hume, Shan Carter, Tom Henighan, and Christopher Olah. Towards
326 monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*,
327 2023. <https://transformer-circuits.pub/2023/monosemantic-features/index.html>.
- 328 [9] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural
329 networks. *arXiv preprint arXiv:1803.03635*, 2018.
- 330 [10] Mengxia Yu, De Wang, Qi Shan, Colorado J Reed, and Alvin Wan. The super weight in large language
331 models. *arXiv preprint arXiv:2411.07191*, 2024.
- 332 [11] Leo Gao, Achyuta Rajaram, Jacob Coxon, Soham V Govande, Bowen Baker, and Dan Mossing. Weight-
333 sparse transformers have interpretable circuits. *arXiv preprint arXiv:2511.13653*, 2025.
- 334 [12] Chris Olah, Nicholas L. Turner, and Tom Conerly. A toy model of interference weights. *Transformer*
335 *Circuits Thread*, July 2025. Informal note.
- 336 [13] Dan Braun, Lucius Bushnaq, Stefan Heimersheim, Jake Mendel, and Lee Sharkey. Interpretability in pa-
337 rameter space: Minimizing mechanistic description length with attribution-based parameter decomposition.
338 *arXiv preprint arXiv:2501.14926*, 2025.
- 339 [14] Jacob Drori. Weight-sparse circuits may be interpretable yet unfaithful. *LessWrong*, February 2026.
340 Accessed: 2026-04-29.
- 341 [15] Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O’Brien, Eric
342 Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia:
343 A suite for analyzing large language models across training and scaling. In *International conference on*
344 *machine learning*, pages 2397–2430. PMLR, 2023.
- 345 [16] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda
346 Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac
347 Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse,
348 Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathe-
349 matical framework for transformer circuits. *Transformer Circuits Thread*, 2021. [https://transformer-
350 circuits.pub/2021/framework/index.html](https://transformer-circuits.pub/2021/framework/index.html).
- 351 [17] Michael T Pearce, Thomas Dooms, Alice Rigg, Jose M Oramas, and Lee Sharkey. Bilinear mlps enable
352 weight-based mechanistic interpretability. *arXiv preprint arXiv:2410.08417*, 2024.

- 353 [18] David Lindner, János Kramár, Sebastian Farquhar, Matthew Rahtz, Tom McGrath, and Vladimir Mikulik.
354 Tracr: Compiled transformers as a laboratory for interpretability. *Advances in Neural Information*
355 *Processing Systems*, 36:37876–37899, 2023.
- 356 [19] Gail Weiss, Yoav Goldberg, and Eran Yahav. Thinking like transformers. In *International Conference on*
357 *Machine Learning*, pages 11080–11090. PMLR, 2021.
- 358 [20] Nelson Elhage, Tristan Hume, Catherine Olsson, Neel Nanda, Tom Henighan, Scott Johnston, Sheer
359 ElShowk, Nicholas Joseph, Nova DasSarma, Ben Mann, Danny Hernandez, Amanda Askell, Kamal
360 Ndousse, Andy Jones, Dawn Drain, Anna Chen, Yuntao Bai, Deep Ganguli, Liane Lovitt, Zac Hatfield-
361 Dodds, Jackson Kernion, Tom Conerly, Shauna Kravec, Stanislav Fort, Saurav Kadavath, Josh Ja-
362 cobsen, Eli Tran-Johnson, Jared Kaplan, Jack Clark, Tom Brown, Sam McCandlish, Dario Amodei,
363 and Christopher Olah. Softmax linear units. *Transformer Circuits Thread*, 2022. [https://transformer-](https://transformer-circuits.pub/2022/solu/index.html)
364 [circuits.pub/2022/solu/index.html](https://transformer-circuits.pub/2022/solu/index.html).
- 365 [21] Lucius Bushnaq, Dan Braun, and Lee Sharkey. Stochastic parameter decomposition. *arXiv preprint*
366 *arXiv:2506.20790*, 2025.
- 367 [22] Brianna Chrisman, Lucius Bushnaq, and Lee Sharkey. Identifying sparsely active circuits through local
368 loss landscape decomposition. *arXiv preprint arXiv:2504.00194*, 2025.
- 369 [23] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang,
370 Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The pile: An 800gb dataset
371 of diverse text for language modeling. *ArXiv*, abs/2101.00027, 2020.
- 372 [24] CodeParrot Team. Codeparrot dataset. [https://huggingface.co/datasets/transformersbook/](https://huggingface.co/datasets/transformersbook/codeparrot)
373 [codeparrot](https://huggingface.co/datasets/transformersbook/codeparrot), 2022.

374 A Appendix

375 A.1 Additional details about corpus selection

376 We use CodeParrot [24] as a proxy for the original Python pretraining data of the Gao et al. (2025)
 377 model, which is not publicly released, and stream The Pile [23] for Pythia-70m, which matches the
 378 pretraining distribution but not the exact training sequences. For the Drori (2026) sparse and dense
 379 models we use SimpleStories [14], which is the actual training corpus.

380 To verify that our chosen corpora are reasonable approximations of each model’s training distribution,
 381 we compare the cross-entropy loss of each model on its evaluation corpus against the loss on
 382 deliberately mismatched corpora. Table A1 reports the comparison. The Gao et al. model achieves
 383 a CE of 2.34 on CodeParrot, but its loss rises sharply on out-of-distribution text. SimpleStories
 384 English yields a CE of 3.70, and French Wikipedia yields 4.50, both substantially higher than the
 385 in-distribution CE. Likewise, the Drori sparse model achieves a CE of 2.98 on SimpleStories, well
 386 below what either model achieves outside its training domain. We take the size of the in-distribution-
 387 versus-out-of-distribution gap as evidence that CodeParrot and SimpleStories are close enough to
 388 each model’s training distribution to serve as evaluation corpora for our pipeline.¹

Table A1: Cross-entropy loss of each model on in-distribution and out-of-distribution corpora. The large gap between in- and out-of-distribution CE indicates that our chosen evaluation corpora are reasonable approximations of each model’s training distribution.

Model	Corpus	CE Loss	In-dist?
Gao et al. (2025)	CodeParrot (Python)	2.34	Yes
Drori (2026), sparse	SimpleStories	2.98	Yes
Gao et al. (2025)	SimpleStories (English)	3.70	No
Gao et al. (2025)	French Wikipedia	4.50	No

389 The per-weight ablation effect $|\Delta\text{CE}|$ spans multiple orders of magnitude across all models, from
 390 $\sim 10^{-3}$ down to below 10^{-6} (Figure 1).



Figure A1: Another representative weight from the Gao et al. (2025) sparse code transformer.

391 A.2 Distributions of recovery and inverse

392 Figures A2 and A3 show the distributions of best-of- N recovery and inverse scores across all
 393 sampled weights, sorted within each model. These curves give a more complete view of the behavior
 394 summarized by the headline interpretability rates. The sparse models have a long tail of weights that
 395 pass both the recovery and inverse criteria, while the dense baselines reach the canonical thresholds
 396 only on a much smaller slice. Negative inverse values for the sparse models reflect predicates whose
 397 negation raises loss above the fully-ablated baseline, which is direct evidence that the predicate is
 398 identifying a real, localized effect rather than partitioning a uniform contribution.

¹Relevant code and representative resources are made available at [anonymous121512.github.io/](https://github.com/anonymous121512)

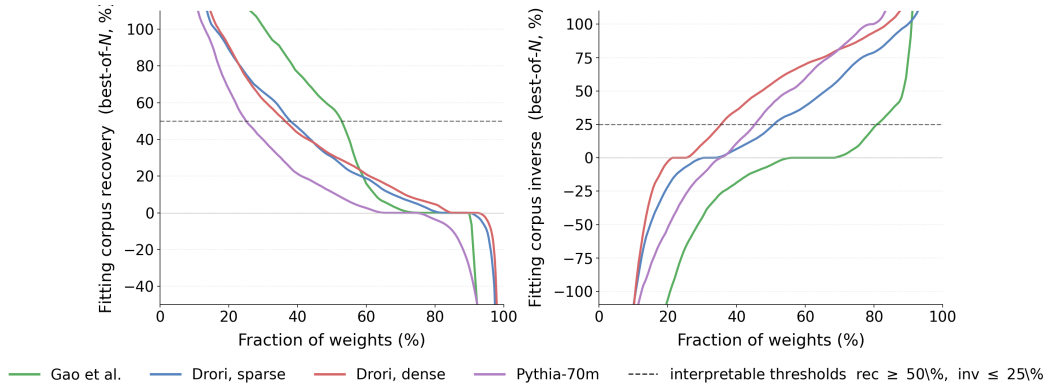


Figure A2: Best-of- N recovery (left) and inverse (right) across all sampled weights, sorted within each model so curves are monotone. Dashed lines mark the canonical thresholds (recovery ≥ 0.5 , inverse ≤ 0.25).

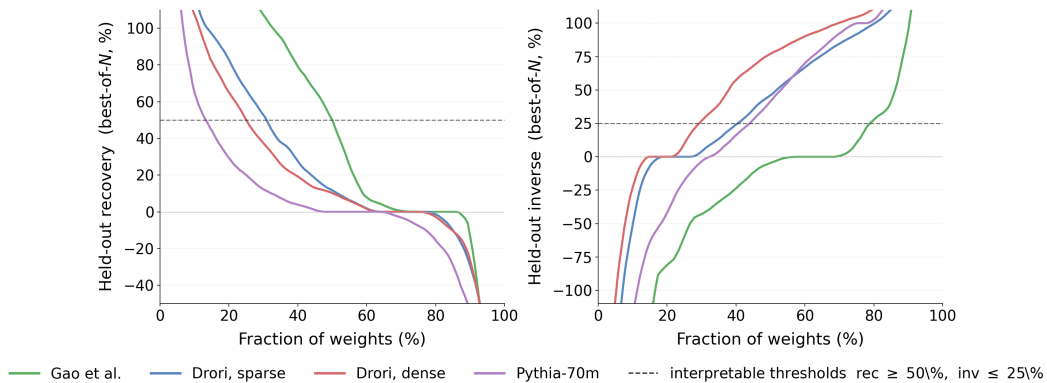


Figure A3: Best-of- N recovery (left) and inverse (right) on the held-out corpus. Same setup as Figure A2, but evaluated on a corpus never seen when generating the candidate predicates.

399 Comparing Figure A2 to the held-out version in Figure A3, the gap between the sparse and dense
 400 models widens. The sparse models retain most of their best-of- N recovery on the held-out corpus,
 401 while the dense models lose more of theirs, which suggests that predicates fitted to sparse-model
 402 weights generalize better than those fitted to dense-model weights. We read this as further evidence
 403 that the predicates for sparse weights are tracking a real property of the weight rather than overfitting
 404 to features of the corpus used to generate them.

405 Figure A4 reports structural statistics of the predicates themselves. Predicates remain short across all
 406 models (median ≤ 4 lines, and only 2 lines on the Drori sparse and dense models), and the predicate
 407 language is the same on sparse and dense models. The difference in interpretability rates is therefore
 408 not driven by sparse models receiving longer or more elaborate predicates.

409 A.3 Marginal shape of s_B vs. per-sample consistency

410 This appendix decomposes the gap and explains why we treat the marginal-shape effect as the
 411 dominant driver.

412 **Per-sample consistency.** Figure A5 fits an ordinary least-squares regression $s_B \sim s_A$ per model
 413 and reports the slope with a 95% CI band. A slope of 1 would indicate that the held-out score
 414 perfectly tracks the in-sample score; a slope of 0 would indicate that s_A has no predictive content for
 415 s_B and the predicate's apparent score on A is accidental. The fitted slopes order in the same direction
 416 as the headline rates: Gao +0.68, Drori sparse +0.50, Drori dense +0.23, Pythia +0.10.

Predicate-length comparison across models

Model	Function length (chars)			Body
	p25	p50	p75	lines
● Gao et al. 2025	143	186	244	4
● Drori 2026 (sparse)	131	162	227	2
● Drori 2026 (dense)	133	164	208	2
● Pythia-70m (dense)	166	229	272	4

What the predicates look like. They fall into a small family of structural patterns. Most are lookup-style tests on a *single* token—“is this token one of {th, ode, b}”, “is it a digit”, “is it a mid-word lowercase fragment”. A smaller set conditions on the immediate neighbour (e.g. “previous word ends in -ed”, “next token is an inflectional suffix”). Even the longer predicates rarely look beyond two or three adjacent tokens: we don’t see explanations like “this token is the subject of the sentence” or “we are inside an open string literal” that would require the model to keep state across the prompt.

Figure A4: Structural statistics of the LLM-generated predicates per model: distribution of predicate length (lines of code) and predicate class (single-token lookup, neighbor test, multi-token context test).

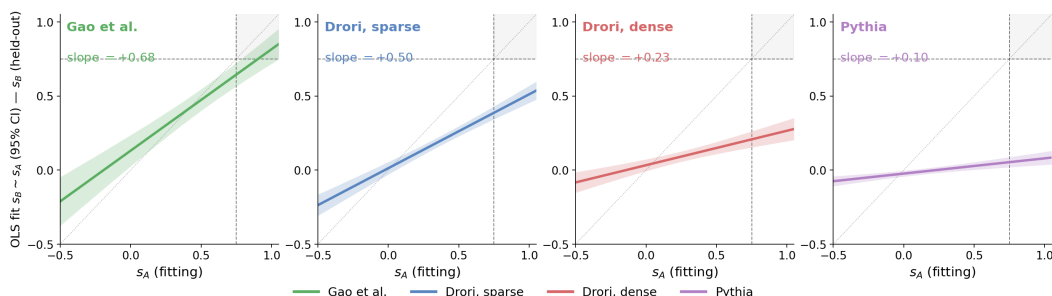


Figure A5: Held-out interpretability score $s_B = \min(\text{rec}_B, 1 - \text{inv}_B)$ across models. OLS fit $s_B \sim s_A$ per model with 95% CI; slope annotated.

417 **Random-chance baseline.** A useful way to see why the marginal effect dominates is to ask what
418 interpretability scores a random predicate would produce. Consider a predicate that fires on a
419 uniformly random subset of corpus positions with coverage p_{fire} , independent of the weight’s actual
420 firing pattern. Restoring the weight at those positions recovers approximately p_{fire} of the ablation
421 effect ($\text{rec} \approx p_{\text{fire}}$); restoring at the complement recovers approximately $1 - p_{\text{fire}}$ ($\text{inv} \approx 1 - p_{\text{fire}}$);
422 and the score is $s \approx \min(p_{\text{fire}}, 1 - p_{\text{fire}}) = p_{\text{fire}}$. Under the gate ($p_{\text{fire}} \leq 0.5$) the random-chance score is
423 bounded above by 0.5, well below the interpretability threshold $s \geq 0.75$. A weight whose predicate
424 scores high on A purely by chance therefore lives near 0.5 on the random-chance axis and re-scoring
425 on B pulls it toward the random-chance mean rather than near 1. Dense-model weights contribute
426 more mass near 0.5 and less mass near 1 on A , so when they are re-scored on B they regress toward
427 the random-chance mean and lose more in the joint criterion than sparse-model weights do. This is
428 why $\text{Interp}_{A \cap B}$ widens the sparse–dense gap relative to Interp_A , and why the widening is mostly
429 attributable to the marginal-shape difference rather than to a per-weight noise difference.

430 A.4 Sensitivity to the candidate budget

431 A single LLM call is too noisy to call a weight interpretable. Pooling N candidates per weight and se-
432 lecting the highest-scoring one tightens rates substantially. Sweeping $N \in \{1, 3, 5, 10, 20, 100, 150\}$
433 (Figure A6) shows interpretability rates rising monotonically with N and saturating before $N = 150$
434 for the sparse models. We use $N = 100$ as the default because the marginal gain per additional
435 candidate flattens past that point and per-weight LLM cost scales linearly. The cross-model ordering
436 is preserved at every N , so the result is not an artifact of the budget we chose.

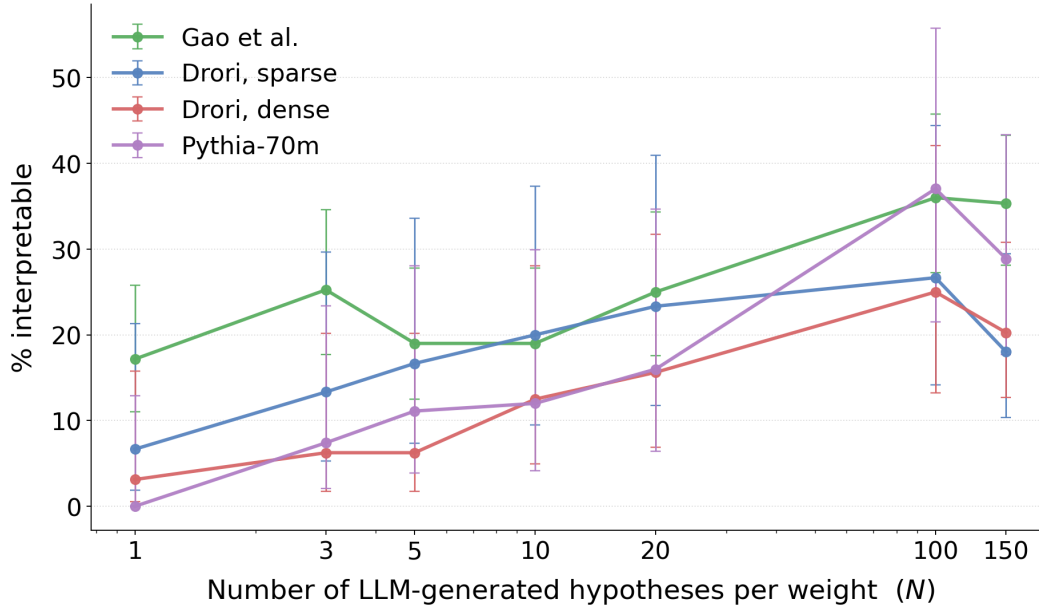


Figure A6: Interpretable rate as a function of the candidate budget N , with 95% Wilson CIs. Single-call N from 1 to 150.

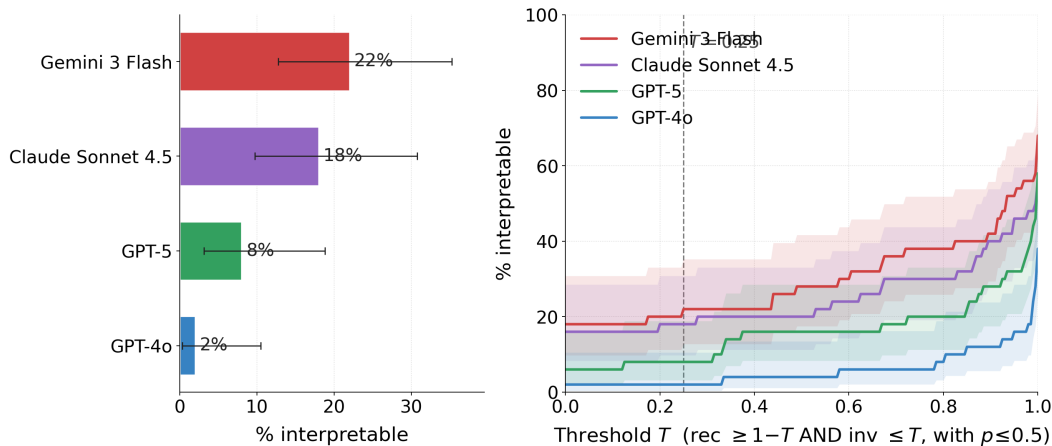


Figure A7: Interpretable rate by auto-interp LLM. **Left:** fraction of sampled weights interpretable at $T = 0.25$, with 95% Wilson CIs, on the same weights evaluated by four different LLMs. **Right:** interpretable rate as a function of threshold T .

437 A.5 Sensitivity to the auto-interp LLM

438 To check that the result is not specific to one auto-interp LLM, we re-run the pipeline with Claude
 439 Sonnet 4.5, GPT-5, and GPT-4o on the same sampled weights.

440 Figure A7 reports the interpretable rate at the canonical threshold $T = 0.25$ for each (model, LLM)
 441 pair, and the rate-vs- T curves at $N = 5$. The shape of the curve is consistent across LLMs. Absolute
 442 rates differ by up to ~ 30 percentage points, with stronger LLMs higher, but the cross-model ordering
 443 (sparse > dense same-architecture > Pythia) is preserved at every T .

444 A.6 Digit-detection neuron: full decomposition

445 This appendix expands the digit-neuron case study of §5. We identify an MLP neuron in the Gao
 446 et al. 2025 sparse code transformer that fires predominantly on numeric tokens; at the activation

447 level the neuron looks monosemantic and its role is naturally read off as “detects digits.” Gao et
 448 al. note that not all nodes in their models are exactly monosemantic on the pretraining distribution
 449 (their Figure 39), and we observe the same on our evaluation corpus — not every node admits a clean
 450 single-concept label at the activation level. We ask what our pipeline returns when applied one level
 451 deeper, to the individual W_{in} entries that connect the neuron to its upstream residual channels.

452 **Per-weight predicates.** The neuron has nine nonzero input weights. Our pipeline assigns a
 453 separately-interpretable predicate to each, and the predicates cluster into three digit-count groups: a
 454 single-digit detector (1 weight), two-digit detectors (5 weights), and a three-digit detector (1 weight);
 455 two further weights receive “generic” or “other” numeric-context predicates that do not split cleanly
 456 by digit count. The strongest weights in each group have high recovery and low inverse, and the digit-
 457 count gates together account for the bulk of the neuron’s output (Figure 6 in §5). The activation-level
 458 “detects digits” label is therefore carrying information about *which* digit count is being detected: the
 459 same neuron’s individual parameters partition the digit-detection role into per-count gates that the
 460 activation-level summary collapses into a single label.

461 **Comparison with the upstream-channel view.** For each of the input weights, we recorded two
 462 things on the same CodeParrot corpus: the tokens at the highest-KL positions when the weight is
 463 ablated, and the tokens at the highest-|activation| positions of the weight’s upstream residual channel.
 464 Table A2 reports the comparison.

Table A2: For each input weight of the digit-detection neuron (neuron 1863), representative top tokens at the weight’s highest-KL positions (left) and top tokens at the upstream channel’s highest-activation positions (right), drawn from the same CodeParrot corpus. Tokens are decoded from the tinypython_2k tokenizer. Rows ordered by max KL.

ch (<i>i</i>)	max KL	weight top tokens	upstream-channel top tokens
464	26.8	100, 200, 300, 111	110, 200, 120, 100
71	22.7	64, 50, 59, 60	<i>non-printable byte</i> , <code>_len</code>
395	17.1	28, 08, 27, 35	39, 28, 35
799	16.1	32, 56, 41, 46	<code>ran</code>
578	13.7	16, 12, 13, 17	<code>tern</code>
131	8.0	4, 2, 3, e	<code>whi</code>
678	6.0	7, 8	<code>!=</code> , <code>raise</code>
87	1.9	90, 89, 95, f	<code>_map</code> , <code>_list</code>

465 For two of the input weights (channels 464 and 395) the upstream channel and the weight pick out
 466 broadly the same tokens (multi-digit numbers in both columns). For the remaining six, the upstream
 467 channel’s top-activating tokens are unrelated to digits (`whi`, `ran`, `tern`, `!=`, `raise`, `_map`, `_list`,
 468 plus one channel dominated by a non-printable byte fragment). The weight’s effect on numeric tokens
 469 does not appear in the channel’s top activations: a description of channel 131 as “fires on `whi`” would
 470 not predict that ablating $w_{131,1863}$ disrupts predictions on the digit 4. We report this as an observation
 471 on a single neuron in a single sparse model; whether it generalizes is left to future work.

472 A.7 String-closing circuit: context-dependent firing

473 A weight participating in the quote-closure circuit of Gao et al. 2025 suppresses the model’s tendency
 474 to predict a closing quote at positions that are still inside an open string. This weight’s importance
 475 depends on context and cannot be inferred from the focus token alone: when we constrain the
 476 predicate language to token-local tests (only the focus token, no look-back), the best predicate the
 477 LLM finds reaches an interpretability score of -0.98 , because whether a `"` token opens or closes a
 478 string is undecidable from the token alone. Allowing a small look-back window (`TOKENS[pos-k]` for
 479 $k \leq 8$) takes the score to $+0.83$, well above the interpretability threshold. The weight is interpretable,
 480 but only under a predicate language rich enough to look back; the same weight under the token-local
 481 language would be filed under “not interpretable,” and the circuit would look incomplete from the
 482 weights up. We use this finding as the empirical justification for the ± 8 -token context window in
 483 §3.2.

484 Figure A8 expands the example to the two quote-detector neurons identified by Gao et al. 2025
 485 — neuron 863 (double quotes) and neuron 2790 (single quotes) — with the assigned predicate,

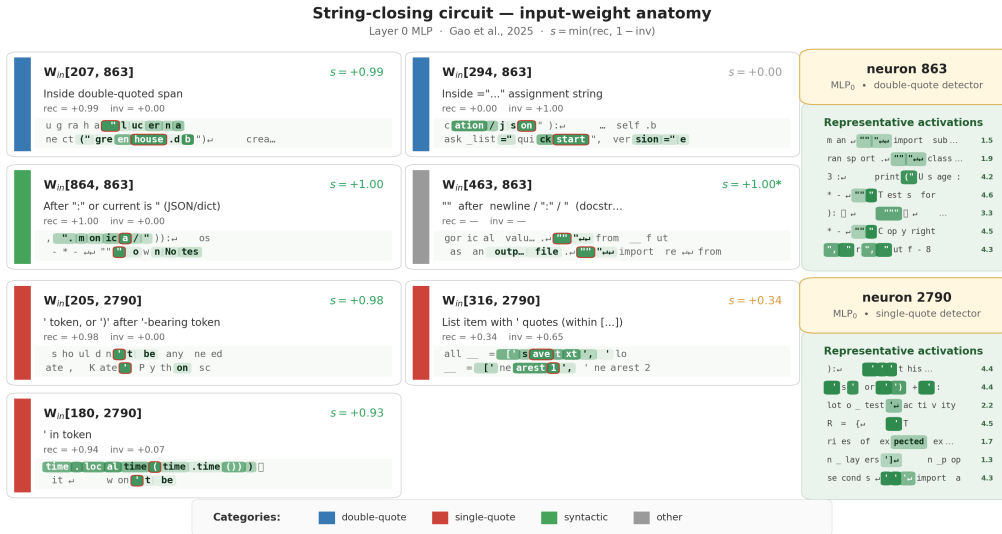


Figure A8: Quote-detector neurons from the string-closing circuit of Gao et al. 2025. **Neuron 863** (top band) fires on double quotes; **neuron 2790** (bottom band) fires on single quotes. Each card shows one nonzero input weight: the predicate, the interpretability score $s = \min(\text{rec}, 1 - \text{inv})$, the recovery and inverse values, and two top-KL example contexts (focus token outlined in red, pill saturation graded by per-token KL within the row). The right panel of each band shows representative activating contexts of the neuron, randomly sampled from the top-90% cumulative-activation band on CodeParrot.

Note: * on the score marks weights whose $|\Delta\text{CE}|$ is at the noise floor, so the rec/inv ratios are not numerically reliable.

486 interpretability score $s = \min(\text{rec}, 1 - \text{inv})$, and two top-KL example contexts for each of their
 487 nonzero input weights, alongside a panel of representative activations of each neuron drawn from the
 488 top-90% cumulative-activation band.

489 A.8 Extended case study: a Drori-sparse speech-verb neuron

490 This appendix complements the case studies of §5 with a fully worked example on the Drori sparse
 491 model: layer-1 MLP pre-GELU neuron 2461. We show the activation-level view first, then decompose
 492 it weight-by-weight, and end by stating what the weight view adds over the activation view.

493 **Activation-level view.** On the SimpleStories corpus, the pre-GELU activation of neuron 2461 has
 494 mean -2.0 , std 1.1 , and a strong skew: the top positive activations sit between $+3.5$ and $+3.9$, and
 495 the top negative activations between -7.2 and -6.4 . Figure A9 shows five context snippets per
 496 direction, randomly sampled from the top-95% activation band; pill saturation tracks magnitude within
 497 each panel. Positive activations (left) are concentrated on speech verbs — *announced*, *whispered*,
 498 *said*, *urged*, *says* — in their natural narrative contexts (“...,” *she announced*, “...” *he whispered*).
 499 Negative activations (right) are dominated by the copulas *is* and *are*. At the activation level, this is a
 500 *speech-verb detector*: it fires up on speech-act verbs and is pushed strongly down on stative-grammar
 501 tokens, in line with what neuron-level autointerp would label.

502 **Weight-level view.** The neuron has 86 nonzero input weights. We rank them by max-KL on a
 503 5,000-sequence SimpleStories sample and profile the top 20, applying the conditional-zero scoring
 504 of §3.3 to each individually; the remaining 66 weights have ablation effects below the noise floor of
 505 our metric on this corpus and are not analyzed here. Figure A10 shows all 20 cards, grouped by the
 506 functional role assigned by their predicate. Five clusters emerge:

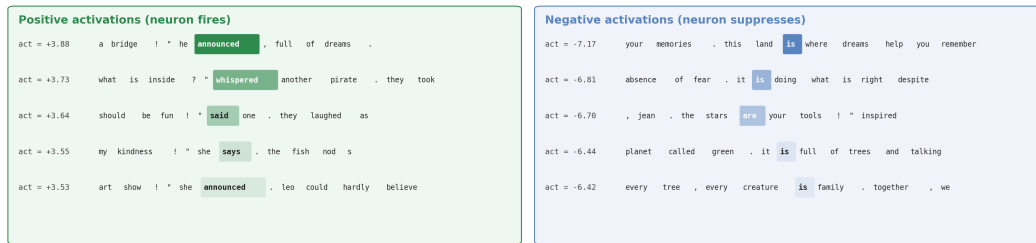


Figure A9: Activation-level view of Drori-sparse layer-1 neuron 2461. **Left:** five contexts where the pre-activation is most positive (random sample from the top-95% positive band). The neuron fires on speech-act verbs. **Right:** five contexts where it is most negative; the neuron is pushed down on copulas (is, are). Pill saturation tracks $|\text{activation}|$ within each panel.

- 507
- 508
- 509
- 510
- 511
- 512
- 513
- 514
- 515
- 516
- 517
- 518
- 519
- 520
- **Speech-verb activators** (positive weights): six channels whose predicates fire on speech verbs — *said, roared, whispered, called, wondered, shouted*. These are the weights that drive the neuron up on the left panel of Figure A9.
 - **Mental-verb suppressors** (negative weights): three channels whose predicates fire on *thought, believed, understood, thanks*. These actively pull the neuron down on a near-neighbour class of verbs that would otherwise share representational support with speech verbs.
 - **Modal-verb suppressor:** one channel firing on *would, could*.
 - **Punctuation gates:** seven channels whose predicates fire on clause-boundary commas and periods. These contribute to the neuron’s negative tail and explain a large share of the is/are suppression observed in the activation view: copulas in this corpus tend to sit one or two tokens after a clause-boundary punctuation mark.
 - **Specific-token gates:** three channels with narrower predicates — a dedicated ‘as’ detector, a ‘called out’ pattern, and a physical-action verb gate.

521 Each card in Figure A10 reports recovery, inverse, and the fraction of corpus positions where the
 522 predicate fires; the strongest cards (*e.g.*, ch. 909 ‘as’ detector with $\text{rec} = 1.22$, $\text{inv} = -0.22$) sit
 523 comfortably above the $T = 0.25$ interpretability bar. Weaker cards (*e.g.*, the modal suppressor on
 524 ch. 38 with $\text{rec} = 0.27$, $\text{inv} = 0.73$) do not pass the bar individually, but their predicates remain
 525 consistent with the role inferred from neighbouring weights in the same group.

526 **What the weight view adds.** At the activation level neuron 2461 reads as a single “speech-verb
 527 detector.” At the weight level the same neuron implements at least *multiple* parallel gating decisions:
 528 it activates speech verbs, suppresses near-neighbour verb classes (mental verbs, modals) that share
 529 representational support with speech verbs, and routes through clause-boundary punctuation patterns
 530 that account for its strongest negative activations. These contributions cancel and add into the single
 531 scalar that activation-level analyses report, so they are not separable from the activation summary
 532 alone — only from the per-weight decomposition.



Figure A10: Weight-level decomposition of Drori-sparse layer-1 neuron 2461. The neuron has 86 nonzero input weights; the figure shows the 20 strongest by max-KL on a 5,000-sequence SimpleStories sample, grouped by the functional role assigned by their predicate. Five context snippets per card are sampled from the top-95% cumulative-KL band of that weight; pill saturation tracks the row's KL within the card. Recovery, inverse, and % active come from the conditional-zero scoring of §3.3.