

# WEAKENING NEURONS: A NEWLY DISCOVERED READ-WRITE FUNCTIONALITY IN TRANSFORMERS WITH OUTSIZE INFLUENCE

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We introduce a new mechanistic interpretability method for gated neurons, based on the cosine similarities between their weight vectors, and use it to gain a number of novel insights into the inner workings of transformer models. First, our method allows us to discover a class of neurons – *weakening* neurons – with surprising behavior: even though there are few, they activate extremely often and have a large influence on model behavior. Second, we show that nine different LLMs have similar patterns with respect to weakening neurons: weakening neurons appear mostly in late layers whereas their counterparts, (*conditional*) *strengthening* neurons, are very frequent in early-middle layers. Third, weakening neurons have a strong effect on model output when gate values are negative – which is surprising since negative gate values are not expected to encode functionality. Thus, for the first time, we observe a mechanism important for transformer functionality that involves negative gate values.<sup>1</sup>

## 1 INTRODUCTION

Mechanistic interpretability research attempts to reverse-engineer the *mechanisms* inside neural networks, such as transformer-based (Vaswani et al., 2017) large language models (LLMs). Some of this work has addressed the interpretation of MLP sublayers, and we follow this line of research. Unlike previous work, we focus on gated activation functions (Shazeer, 2020), which are used in recent LLMs like OLMo, Llama and Gemma, but so far lack an extensive analysis from the interpretability perspective.

Much previous work analyzes neurons<sup>2</sup> based only on the contexts in which they activate (Voita et al., 2024) or based only on their output weights<sup>3</sup> (Gurnee et al., 2024). However, neither of these fully captures the *mechanisms* that neurons implement: we also have to understand the relationship between input and output behavior of neurons, what we call their read-write (RW) functionality. We put this in the center of our analysis, and propose a simple method to investigate RW functionalities: computing cosine similarities between input (reading) and output (writing) weights.

This new approach allows us to gain a number of striking novel insights into the inner workings of LLMs. In particular, we discover a small class of neurons – *weakening* neurons – with outsize influence and often surprising behavior.

Our contributions are as follows: (i) We are the first to investigate read-write behavior of gated neurons, using cosine similarities of weight vectors. (ii) Applying this method to nine LLMs, we observe universal patterns: Early-middle layers contain many *conditional strengthening* neurons, and late layers tend more towards *weakening*. (iii) Thanks to the RW perspective, we discover a small class of neurons (*weakening* neurons), that is highly influential in often surprising ways: they activate often (in the sense of having a gate value above zero), and they influence various metrics, even in earlier layers where they are very rare. (iv) We introduce a new method of conditional ablation that enables us to find *which activations* of a given neuron are responsible for a certain behavior. (v)

<sup>1</sup>We publish code at [https://anonymous.4open.science/r/RW\\_functionalities-4D32](https://anonymous.4open.science/r/RW_functionalities-4D32).

<sup>2</sup>We use "neuron" to refer to a hidden dimension inside the MLP layer.

<sup>3</sup>We use "weight" to refer to a weight vector, not a scalar.

Applying this method to weakening neurons, we find that some of their effect is due to cases in which their gate value is negative. Thus, for the first time, we observe a mechanism involving negative values of the Swish activation function.

## 2 RELATED WORK

There is a large body of work on interpretability of transformer-based LLMs. Elhage et al. (2021) introduce the notion of residual stream. nostalgebraist (2020), Belrose et al. (2023) propose to interpret residual stream states as intermediate guesses about the next token; Rushing & Nanda (2024) discuss this as the *iterative inference hypothesis*. On a similar note, many works hypothesize that directions in model space can correspond to concepts; Park et al. (2024) discuss this as the *linear representation hypothesis*. Lad et al. (2024) define *stages of inference*. Similar to our work, Elhelo & Geva (2024) investigate input-output functionality of heads (instead of neurons).

Much research has attempted to understand individual neurons. Geva et al. (2021) present them as a key-value memory. Other neuron analysis work includes (Miller & Neo, 2023; Niu et al., 2024).

The focus on individual neurons has been criticized. Morcos et al. (2018) find that in good models, neurons are not monosemantic (but for image models, not LLMs). Millidge & Black (2022) find interpretable directions that do not correspond to individual neurons. Elhage et al. (2022) argue that interpretable features are non-orthogonal directions in model space and can be superposed. This corresponds to sparse linear combinations of neurons in MLP space. This has inspired a series of work on sparse autoencoders (SAEs), starting with Sharkey et al. (2022).

The focus on SAEs has been criticized: recent studies indicate that they do not always outperform baselines (Kantamneni et al., 2025; Leask et al., 2025; Mueller et al., 2025; Wu et al., 2025). A middle ground is possible: Gurnee et al. (2023) argue that interpretable features correspond to sparse combinations of neurons; this includes 1-sparse combinations, i.e., individual neurons. Accordingly, there is recent work on new classes of interpretable neurons, e.g. Ali et al. (2025); Zhao et al. (2025).

Several works classify neurons based on the **contexts** in which they activate (Voita et al., 2024; Gurnee et al., 2024). For example, Voita et al. (2024) find *token detectors* that suppress repetitions. Gurnee et al. (2024) also define *functional roles* of neurons based on their **output** weight vector, such as *suppression neurons* that suppress a specific set of tokens. Stolfo et al. (2024) also investigate some output-based neuron classes.

There has been less focus on the input-output perspective. Gurnee et al. (2024) compute cosine similarities between input and output weights for GPT-2 (Radford et al., 2019), but do not interpret their results. Elhage et al. (2021) mention the idea of input-output analysis (footnote 7), but do not follow up. Note that input-output analysis for gated activation functions adds complexity because, in addition to input and output weight vectors, the gating mechanism is crucial for RW functionality.

## 3 PRELIMINARIES

### 3.1 GATED ACTIVATION FUNCTIONS

We work on *gated activation functions* like SwiGLU or GEGLU (Shazeer, 2020). Gated activation functions are used widely, e.g., OLMo (Groeneveld et al., 2024) and Llama (Touvron et al., 2023) use SwiGLU, and Gemma (Gemma, 2024) uses GEGLU. Here we briefly describe SwiGLU. GEGLU replaces Swish with GELU, but is otherwise identical.

Traditional activation functions like ReLU require one weight matrix on the input side and one on the output side: The MLP outputs

$$\mathbf{W}_{\text{out}} \text{ReLU}(\mathbf{W}_{\text{in}} \mathbf{x}_{\text{norm}}),$$

where ReLU is applied element-wise to each neuron (it takes a single scalar as argument).<sup>4</sup>

<sup>4</sup>We write  $\mathbf{x}_{\text{norm}}$  for the residual stream state before the MLP layer of interest (with layer normalization already applied, for the many models that use pre-norm).

Other traditional activation functions are Swish( $x$ ) :=  $x/(1 + \exp(-x))$  (Ramachandran et al., 2017) and GELU( $x$ ) :=  $x\Phi(x)^5$  (Hendrycks & Gimpel, 2016). Both of these can be seen as smooth approximations of ReLU. They are known to work better than ReLU, which is widely believed to be because of their good differentiability (e.g. Lee, 2023), i.e. better training dynamics.

In contrast to these traditional functions, a *gated activation function* like SwiGLU requires two weight matrices on the input side: The MLP outputs

$$\mathbf{W}_{\text{out}} (\text{Swish}(\mathbf{W}_{\text{gate}}\mathbf{x}_{\text{norm}}) \odot (\mathbf{W}_{\text{in}}\mathbf{x}_{\text{norm}})), \quad (1)$$

where  $\odot$  denotes element-wise multiplication (a.k.a. Hadamard product).<sup>6</sup>

We find it more intuitive to separately consider each neuron: The neuron adds the vector

$$\text{Swish}(\langle \mathbf{w}_{\text{gate}}, \mathbf{x}_{\text{norm}} \rangle) \cdot \langle \mathbf{w}_{\text{in}}, \mathbf{x}_{\text{norm}} \rangle \cdot \mathbf{w}_{\text{out}} \quad (2)$$

to the residual stream. Here  $\mathbf{w}_{\text{gate}}$ ,  $\mathbf{w}_{\text{in}}$  are each one of the  $d_{\text{MLP}}$  rows of  $\mathbf{W}_{\text{gate}}$ ,  $\mathbf{W}_{\text{out}}$ , and  $\mathbf{w}_{\text{out}}$  is one of the  $d_{\text{MLP}}$  columns of  $\mathbf{W}_{\text{out}}$ .<sup>7</sup> These weight vectors, as well as  $\mathbf{x}_{\text{norm}}$ , all have dimensionality  $d_{\text{model}}$ .<sup>8</sup>

In this framework, SwiGLU can be described as a function of two scalars:

$$\text{SwiGLU}(x_{\text{gate}}, x_{\text{in}}) := \text{Swish}(x_{\text{gate}}) \cdot x_{\text{in}},$$

Unlike ReLU, gated activation functions can output arbitrary positive or negative values. For example, if  $x_{\text{gate}} > 0$  and  $x_{\text{in}} \ll 0$ , then  $\text{SwiGLU}(x_{\text{gate}}, x_{\text{in}}) \ll 0$ .

### 3.2 WEIGHT PREPROCESSING

Our code uses TransformerLens (Nanda & Bloom, 2022). When a model is loaded, certain preprocessing steps are applied to the weights, in order to make the weights more interpretable without changing model behavior (see their documentation for details).

We propose an additional preprocessing step specific to gated activation functions: For each neuron, we multiply  $\mathbf{w}_{\text{in}}$  and  $\mathbf{w}_{\text{out}}$  by the sign of  $\cos(\mathbf{w}_{\text{gate}}, \mathbf{w}_{\text{in}})$ . See section C for our argument on why we do this and why it does not change model behavior.

## 4 METHOD

### 4.1 APPROACH

We follow the *residual stream* perspective of Elhage et al. (2021): Individual model units *read* from the residual stream and then update it by *writing* (adding) to it. In the case of an MLP neuron, the scalar products  $\langle \mathbf{w}_{\text{gate}}, \mathbf{x}_{\text{norm}} \rangle$  and  $\langle \mathbf{w}_{\text{in}}, \mathbf{x}_{\text{norm}} \rangle$  can be thought of as *reading* how much the residual stream  $\mathbf{x}_{\text{norm}}$  conforms to the *directions*  $\mathbf{w}_{\text{gate}}$  and  $\mathbf{w}_{\text{in}}$ . The neuron then *writes* a multiple of the direction  $\mathbf{w}_{\text{out}}$  to the residual stream.

A semantic interpretation is that a neuron detects a *concept* in the residual stream, and in turn also writes a concept. This semantic interpretation is not a necessary assumption for our neuron classification, but is helpful for building intuition and interpreting results.

This framework leads to our main research question: **What is the relationship between what a neuron reads and what it writes?** We take an approach based on **weights** (as opposed to activations) of **neurons** (as opposed to SAE-like features).

**Weight-based.** There are many ways to address the question; we choose a purely weight-based approach: computing the **cosine similarity of input and output weights**. This lets us understand the mathematical function that a neuron implements in terms of updates to the residual stream.

<sup>5</sup> $\Phi$  is the cumulative distribution function (cdf) of a standard normal distribution.

<sup>6</sup>In other works,  $\mathbf{W}_{\text{in}}$  and  $\mathbf{W}_{\text{out}}$  are also called  $\mathbf{W}_{\text{up}}$  and  $\mathbf{W}_{\text{down}}$ , respectively.

<sup>7</sup>This is assuming the right-to-left notation of our equation (1). With left-to-right notation rows and columns are switched.

<sup>8</sup>Following TransformerLens (Nanda & Bloom, 2022), we write  $d_{\text{model}}$  for the model dimensionality (i.e. the dimensionality of the residual stream), and  $d_{\text{MLP}}$  for the hidden dimensionality of a given MLP layer.

Table 1: Six prototypical read-write (RW) functionalities. See section 4.2 for details.

$\cos(\mathbf{w}_{in}, \mathbf{w}_{out})$	$ \cos(\mathbf{w}_{gate}, \mathbf{w}_{out}) $	$\approx 1$ (or $> 0.5$ )	$\approx 0$ (or $< 0.5$ )
$\approx +1$ (or $> +0.5$ )		strengthening	conditional strengthening
$\approx -1$ (or $< -0.5$ )		weakening	conditional weakening
$\approx 0$ (or $\in [-0.5, +0.5]$ )		proportional change	orthogonal output

**Neuron-based.** This cosine similarity method could in principle also be applied to SAE-style features instead of neurons, as long as each feature has well-defined input and output weights. However, for this paper we decided to just investigate neurons, and defer a possible investigation of SAEs to future work. We do so for the following reasons: (i) As argued in section 2, individual neurons can still be a promising research direction today. (ii) For any given LLM, neurons are readily available and clearly defined, whereas researchers may have published several SAEs with different sizes and architectures – or, for less popular models, no SAE at all. (iii) We expect that findings from neurons will, to some extent, carry over to linear combinations of neurons. See section D for more on this.

In section 5 we will see that, despite being "only" weight-based and neuron-based, our method already yields striking results.

#### 4.2 TAXONOMY OF RW FUNCTIONALITIES

We now think through what different combinations of weight cosine similarities would mean for neuron RW functionality, and introduce our terminology. For the moment we focus on the prototypical cases, in which cosine similarities are approximately  $\pm 1$  or 0. We present a taxonomy of these prototypical RW functionalities in table 1.

Generally, when the output weight is similar enough to (one of) the detected directions, we speak of *input manipulation*, as opposed to **orthogonal output** neurons which write to directions not detected in the input. Intuitively, input manipulator neurons *manipulate* the concept that they detect.

As special cases of input manipulation, we define: (i) **Strengthening** and **weakening** neurons: all three weight vectors are roughly collinear, and specifically  $\cos(\mathbf{w}_{in}, \mathbf{w}_{out}) \approx \pm 1$ . The neuron detects a direction and then adds it to / removes it from the residual stream. (ii) **Conditional strengthening / weakening** neurons:  $\mathbf{w}_{in}$  and  $\mathbf{w}_{out}$  are roughly collinear and  $\mathbf{w}_{gate}$  is orthogonal to them. The neuron also strengthens / weakens the direction detected by its  $\mathbf{w}_{in}$  vector, but will only activate *conditional on  $\mathbf{w}_{gate}$  being present in the residual stream*. (iii) **Proportional change** neurons:  $\mathbf{w}_{out}$  is collinear to  $\mathbf{w}_{gate}$ , but is orthogonal to  $\mathbf{w}_{in}$ . If  $\mathbf{w}_{gate}$  is present in the residual stream, then the neuron writes a *positive or negative* multiple of this direction to the residual stream. This multiple is proportional to the presence of  $\mathbf{w}_{in}$  in the residual stream.

These prototypical classes are limited in scope: Many cosines will not be close to 0 or  $\pm 1$ . For this general case, this paper explores three options to understand neuron RW functionalities at different levels of granularity: (1) Classify neurons according to the closest prototypical case (we choose a threshold  $\tau = \pm 0.5$ ). (2) Plot the marginal distributions of the three cosine similarities. (3) Place neurons in a scatter plot, based on their three weight cosines.

In option 1 (threshold-based classification),  $\cos(\mathbf{w}_{in}, \mathbf{w}_{gate})$  may not always “match” the other two cosine similarities. Consider for example the case of strengthening: In the prototypical case with exact equalities ( $\cos(\mathbf{w}_{in}, \mathbf{w}_{out}) = \cos(\mathbf{w}_{gate}, \mathbf{w}_{out}) = 1$ ), all three weight vectors are collinear, so we also have  $\cos(\mathbf{w}_{in}, \mathbf{w}_{gate}) = 1$ . But without exact equalities (if we just know  $\cos(\mathbf{w}_{in}, \mathbf{w}_{out})$  and  $\cos(\mathbf{w}_{gate}, \mathbf{w}_{out})$  are both above 0.5), it does not follow that  $\cos(\mathbf{w}_{in}, \mathbf{w}_{gate})$  is also above 0.5.<sup>9</sup> When such a “mismatch” occurs, we prepend *atypical* to the category’s name: In this example, we will

<sup>9</sup>For example, the two reading weights may be orthogonal ( $\cos(\mathbf{w}_{in}, \mathbf{w}_{gate}) = 0 < 0.5$ ), but  $\mathbf{w}_{out} = \mathbf{w}_{gate} + \mathbf{w}_{in}$ ; then  $\cos(\mathbf{w}_{in}, \mathbf{w}_{out}) = \cos(\mathbf{w}_{gate}, \mathbf{w}_{out}) \approx 0.71 > 0.5$ .

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

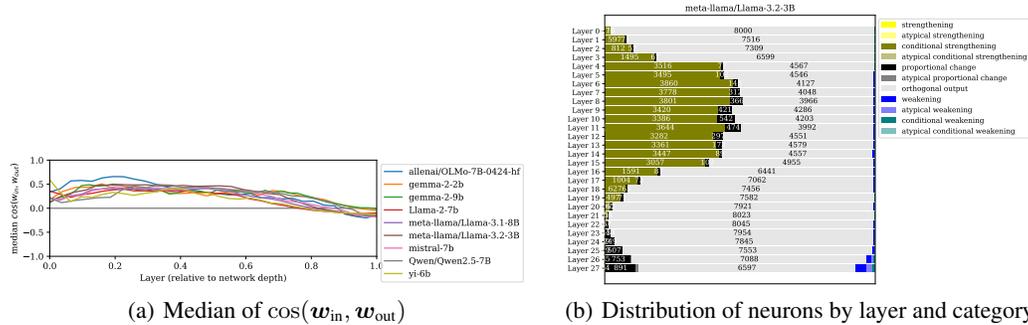


Figure 1: (a) Median of  $\cos(\mathbf{w}_{in}, \mathbf{w}_{out})$  by layer (x-axis) for 9 models of 2B to 9B parameters. For all models, the value is positive in the beginning and negative in the end, indicating that early-middle layers “strengthen” directions they find in the residual stream whereas later layers tend more towards “weakening” them. (b) Distribution of neurons by layer and category.

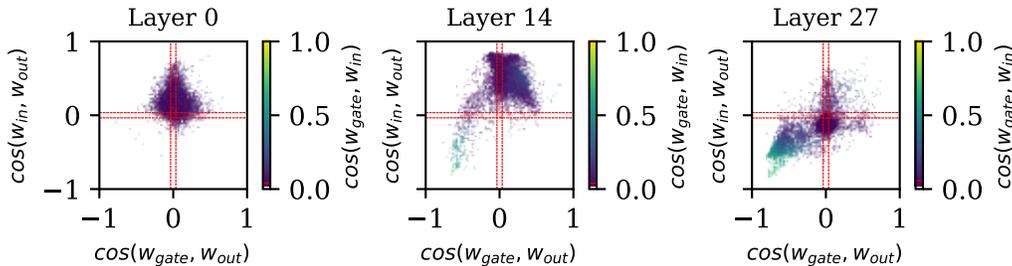


Figure 2: Fine-grained analysis of neuron RW behavior in three layers of Llama-3.2-3B, based on the configuration of their three weight vectors in parameter space. Each subplot represents a layer, each dot a neuron. The red lines mark the 95% randomness regions for each of the three cosine values. (There is a dotted line for variant (i) and a dashed line for variant (ii) in section 4.3, but they are almost the same.)

We see that many neurons are outside the randomness regions, indicating that they manipulate their input in some way. Purple dots at the top of the plots are conditional strengthening neurons. Lighter dots in the bottom left corner are weakening neurons.

speak of an atypical strengthening neuron. In figure 1(b) we will see that such neurons exist, but are quite rare overall.

### 4.3 RANDOM BASELINES

Given a cosine similarity between weight vectors, we test if it is significantly different from random. To do so, we consider two random baselines: (i) i.i.d. Gaussian vector entries (as in a randomly initialized model); (ii) a layer-specific baseline based on “mismatched cosines”. See section E for details. In practice, we find that both baselines give quite similar 95% randomness ranges.

## 5 WHERE TO FIND WEAKENING NEURONS

In this section we compute cosine similarities of neuron weights as described in section 4, to investigate which RW functionalities actually appear in LLMs, and in which layers. Strikingly, our results are **consistent across models**. Across models, we find that there is a small but (as we will see later) influential number of **weakening neurons**, mostly in late layers. Other RW functionalities appear in other ranges: in particular, early-middle layers of all models contain a lot of conditional strengthening neurons.

270 Concretely, we apply our method to 12 LLMs: Gemma-2-2B, Gemma-2-9B (Gemma, 2024), Llama-  
 271 2-7B, -3.1-8B, -3.2-1B, -3.2-3B (Touvron et al., 2023), OLMo-1B, OLMo-7B-0424 (Groeneveld  
 272 et al., 2024), Mistral-7B (Jiang et al., 2023), Qwen2.5-0.5B, Qwen2.5-7B (Yang et al., 2024), Yi-6B  
 273 (01.AI et al., 2025). These models use SwiGLU, except for Gemma, which uses GeGLU.

274 To demonstrate our finding in more detail, we present three representative plots. (See section J for  
 275 more.) Figure 1(a) shows the median value of  $\cos(\mathbf{w}_{in}, \mathbf{w}_{out})$  across all layers of the nine larger  
 276 models. The common pattern is clearly visible: In early-middle layers of all models, a majority of  
 277 neurons has a  $\cos(\mathbf{w}_{in}, \mathbf{w}_{out})$  high above zero, indicating strengthening; in late layers, this median  
 278 cosine similarity goes slightly below zero, indicating a relative majority of weakening neurons.

279 The other two plots focus on **Llama-3.2-3B**, but the patterns we describe are general: see section J  
 280 for other models. Figure 1(b) shows RW class distribution across layers. In figure 2, we plot  
 281 the distribution of neurons in a few selected layers, by displaying each neuron as a point with  
 282  $\cos(\mathbf{w}_{gate}, \mathbf{w}_{out})$  indicated on the x-axis,  $\cos(\mathbf{w}_{in}, \mathbf{w}_{out})$  on the y-axis and  $\cos(\mathbf{w}_{gate}, \mathbf{w}_{in})$  as its color.

283 **Input manipulation.** First, we see that a large proportion of neurons are input manipulators (i.e.,  
 284 they are not orthogonal output neurons): In figure 1(b), these are 25% of all neurons, and as much  
 285 as 50% in early-middle layers (layers 7–11<sup>10</sup>). What is more, figure 2 shows that even the neurons  
 286 classified as "orthogonal output" often belong to clusters that are centered above/below the horizontal  
 287 line. Their weight cosine similarities often exceed the significance threshold. E.g., in layer 14, there  
 288 are many neurons whose  $\cos(\mathbf{w}_{in}, \mathbf{w}_{out})$  (y-axis) is below 0.5 but above the significance threshold.  
 289 This suggests that even the "orthogonal output" neurons perform input manipulation to some extent.

290 **Different RW functionalities.** Weakening neurons represent a large share of the (relatively few)  
 291 input manipulators in late layers. They form a somewhat separate cluster in figure 2 (in the bottom-  
 292 left corner of the rightmost subplot). Another important input manipulator class in late layers is  
 293 proportional change. In contrast, across all models, early middle layers are dominated by conditional  
 294 strengthening. In fact, the majority of input manipulators (more than 80% in Llama) belong to just  
 295 this one class.

296 This general pattern of strengthening-then-weakening holds across models, as figure 1(a) shows at  
 297 one glance. In figure 2 (and figure 52 in the appendix), the pattern manifests as a large cluster of  
 298 neurons, centered clearly above the x-axis in most layers, but moving below it in the last few layers.

299 In summary, we find across models that conditional strengthening dominates in early-middle layers,  
 300 but in late layers we find more weakening neurons.

## 303 6 ABLATION EXPERIMENTS

304 Since model training produced so many input manipulator neurons, we hypothesize that they must  
 305 contribute to model performance in an important way. We now test this hypothesis by ablating  
 306 neurons based on their RW functionality. We find that weakening neurons have the highest effect on  
 307 the metrics that we tested – this is completely unexpected since weakening neurons are a very small  
 308 class of a few hundred neurons.

309 In the rest of the paper, we run a model on a dataset – whereas in section 5 we just applied our  
 310 weight-based method from section 4. Therefore, to save resources, we focus on a single model: We  
 311 choose **OLMo-7B**, because its training dataset, Dolma (Soldaini et al., 2024), is publicly available  
 312 and its RW functionalities mostly follow the typical patterns. As a dataset, we use a random subset of  
 313 20M tokens from Dolma,<sup>11</sup> except for attribute rate, where we follow the setup of Geva et al. (2023).

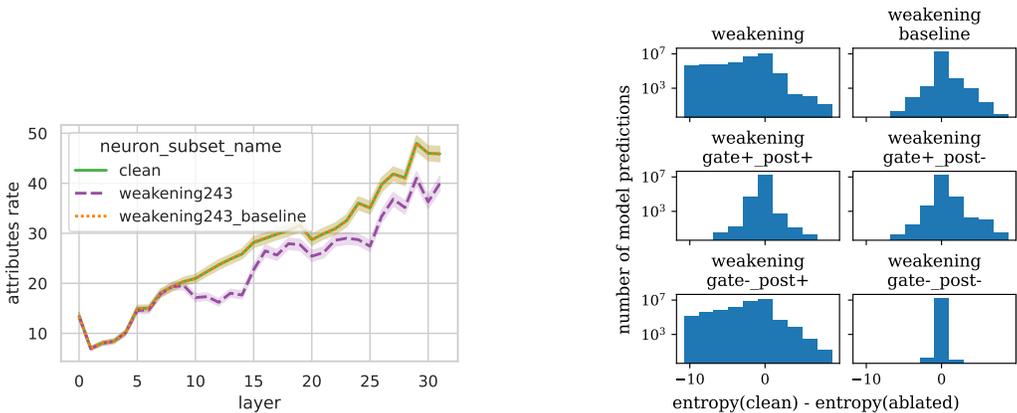
### 316 6.1 EXPLORING RW CLASSES AND METRICS

317 We run the model on our dataset and record various metrics, such as the loss. In each run we ablate a  
 318 number of neurons from a different RW class, or (as a baseline) the same number of *random* neurons  
 319 from the same layers. This enables us to observe the effect of various RW classes on these metrics.  
 320 The baseline verifies if effects are due to the layers rather than RW classes.

322 <sup>10</sup>We use zero-based (Python-style) indexing throughout the paper.

323 <sup>11</sup>The size of 20M tokens follows Voita et al., 2024.

324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377



(a) Effect on attribute rate of ablating all 243 weakening neurons (weakening243), or 243 random neurons from the same layers (weakening243\_baseline). These baseline neurons do not have any sizable influence. In contrast, the effect of weakening neurons is clearly visible, already from layer  $\approx 10$  onward, even though weakening neurons are few and mostly in late layers. In the appendix (figures 14 to 16) we show results for other neuron classes, all of which are indistinguishable from the "clean" line.

(b) Effect on entropy of various neuron activations. E.g., in  $\approx 10^5$  next-token predictions, weakening neurons decrease the entropy by about 10 nats, whereas they increase it much more rarely. "Weakening baseline" denotes random neurons from the same layers as weakening neurons. The bottom four plots describe the results of conditional ablations (section 6.2). E.g., "gate+\_post+" describes the effect of those activations in which  $x_{gate} > 0$  and  $x_{post} > 0$ .

Figure 3: Effect of zero-ablating weakening neurons.

The main metrics we consider are attribute rate (Geva et al., 2023) and entropy of the output distribution. We justify these choices in section F. We try two types of ablation: zero ablation (setting activations to zero), and mean ablation (setting them to the mean activation of the given neuron).

We find that **ablating weakening neurons has a large effect** on both metrics, and this effect is not seen with other classes or with other neurons from the same layers. The effect is clearest with zero ablation, but also present with mean ablation (see section F.4 for mean ablation results). For attribute rate, the effect is most visible in layers  $\approx 10$  and onward. See figure 3(a). This is particularly interesting since there are very few weakening neurons in these early-middle layers. The case of entropy is also striking: Figure 3(b) shows that ablating weakening neurons often makes the output distribution flatter; in other words weakening neurons make the output distribution *sharper*. We would expect the opposite: removing information from the residual stream should make it less informative and therefore flatten the output distribution.

## 6.2 CONDITIONAL ABLATIONS

We now further investigate the effect of weakening neurons on entropy. We use **conditional ablations**: We ablate only some activations of each neuron, based on the signs of the corresponding  $x_{gate}$  and  $x_{in}$ . Specifically, we consider the following four conditions (using the preprocessing from section 3.2): (i)  $x_{gate} > 0, x_{in} > 0$ , leading to  $x_{post} > 0$ ; (ii)  $x_{gate} > 0, x_{in} < 0$ , leading to  $x_{post} < 0$ ; (iii)  $x_{gate} < 0, x_{in} < 0$ , leading to  $x_{post} > 0$ ; (iv)  $x_{gate} < 0, x_{in} > 0$ , leading to  $x_{post} < 0$ .

We find that a large part of the sharpening effect of weakening neurons is due to case (iii): In figure 3(b), case (iii) (bottom left subplot) shows entropy effects similar to those of weakening neurons as a whole, whereas this is much less the case for the other subplots. This is surprising, but also solves the mystery we encountered earlier (i.e., we expected weakening neurons to flatten the distribution, but in reality they often sharpen it).

It is *surprising* for two reasons: First, these negative  $x_{gate}$  activations are relatively rare in weakening neurons (as we will see in section 7). Second, because of the Swish function, **negative gate values** are relatively small (whereas positive values can be arbitrarily large), and it was often assumed they were

only useful for training dynamics (see section 3.1). Our results show for the first time (concurrently with Kong et al. (2025) who focus on a different phenomenon) that negative gate values have a strong effect on model mechanisms (not just training). *This shows that, for mechanistic interpretability research, Swish is not reducible to ReLU.*

This is also *explanatory*: When  $x_{\text{gate}} < 0$ , the usual neuron behavior gets a minus sign in front, so that weakening neurons take on a strengthening behavior. E.g., if a neuron usually detects "minus again" ( $w_{\text{gate}}$ ) and writes "again" ( $w_{\text{out}}$ ), in case (iii) it detects "again" ( $-w_{\text{gate}}$ ) and writes "again" ( $w_{\text{out}}$ ), which indeed makes the output distribution sharper. We will analyze such a neuron in section 8.

### 6.3 CASE STUDY OF ENTROPY REDUCTION

To understand this phenomenon further, we study a particular text example, namely where the entropy reduction by case (iii) activations of weakening neurons was most extreme (with zero ablation).

The input text is: *Yesterday (21 December) the Government announced a package of support for hospitality and leisure businesses that are losing trade because of the O* and the correct next token is *mic* (as in *Omicron*). The model predicts this next token correctly.

Which tokens have the largest score difference between clean and ablated runs? We find that, in the clean run, *mic* and similar tokens get a massive boost (of up to 12 points) compared to the ablated run, whereas no token gets its score reduced by nearly as much. Thus, at least in this case, the case (iii) activations of weakening neurons sharpen the output distribution by boosting the correct next token.

We further investigate whether any single weakening neuron has a  $w_{\text{out}}$  similar to *mic*. This is not the case. This suggests that in this case weakening neurons work together (in superposition, cf. Elhage et al., 2022) to achieve the observed effect.

## 7 WEAKENING NEURONS ACTIVATE OFTEN

Our findings from section 6 raise the question of how often weakening neurons activate, i.e., how often their gate value is positive. In fact, Gurnee et al. (2024) found a negative correlation between activation frequency and  $\cos(w_{\text{in}}, w_{\text{out}})$  – but in a GELU model. We now investigate whether a similar phenomenon occurs with gated activation functions. We show the most striking result in figure 4. We show other results in section J and discuss some of them in section G.

Consistent with Gurnee et al. (2024), we find that the many (conditional) strengthening neurons activate very rarely, and (conditional) **weakening neurons activate very often**. In fact, in most layers there is an almost linear negative relationship between  $\cos(w_{\text{in}}, w_{\text{out}})$  and activation frequency: correlations are at least  $-0.71$  in all layers except the last two (which have  $-0.29$  and  $+0.29$ ).

It seems that each conditional strengthening neuron is responsible only for a narrow domain, perhaps a specific set of tokens. This result is another indication that weakening neurons have a disproportionately large influence on model behavior. Note however that activation frequencies do not fully explain their effect, since we found that even their negative gate values are influential (section 6).

## 8 CASE STUDY OF A WEAKENING NEURON

We now qualitatively examine two neurons in more detail: a strengthening and a weakening neuron. We will see that weakening neurons can have a quite complex behavior. In section I we detail the

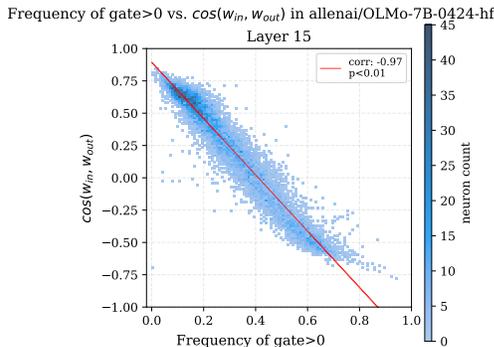


Figure 4: Two-dimensional histogram of activation frequency (x-axis) vs.  $\cos(w_{\text{in}}, w_{\text{out}})$  (y-axis).

432 methods (including how we chose the neurons), present this analysis in more detail, and present many  
433 more case studies from various RW functionalities.

434 To analyze the neurons, we combine the RW perspective with two well-established neuron analysis  
435 methods: projecting weights to vocabulary space (nostalgebraist, 2020; Geva et al., 2022; Dar et al.,  
436 2023; Gurnee et al., 2024; Voita et al., 2024), and finding text examples which strongly activate the  
437 neuron (Dalvi et al., 2019; Geva et al., 2021; Nanda, 2022; Voita et al., 2024; Gurnee et al., 2024).

438 We choose neuron **28.4737** for strengthening and **31.9634** for weakening.<sup>12</sup> Both of them are also a  
439 *prediction* neuron in the sense of Gurnee et al. (2024), which indicates that they directly promote a  
440 specific set of tokens and are thus likely to be monosemantic.

441 We can see that **strengthening neuron 28.4737** has a straightforward read-write behavior: It further  
442 promotes *review* when the residual stream already indicates that this is the obvious next token.

443 In contrast, **weakening neuron 31.9634** is much harder to interpret: The weights indicate that this  
444 neuron produces "*again*" when the residual stream contains "minus *again*"; but the examples strongly  
445 activating the neuron do not have an obvious semantic relationship to *again*.

446 The most interpretable activations were some weaker positive activations when *again* is a plausible  
447 continuation, e.g., on the token **once** (as in *once again*). These are cases with negative  $x_{\text{gate}}$  values  
448 (and also  $x_{\text{in}} < 0$ , hence positive activations) – a case that we found to be important in section 6.2.  
449 In these cases, *again* is already weakly present in the residual stream before the last MLP, and the  
450 neuron reinforces *again*. Thus the behavior of this particular weakening neuron is interpretable in the  
451  $x_{\text{gate}} < 0$  case, echoing our finding from section 6.2 that this case is surprisingly relevant to model  
452 behavior.

453 These two case studies show that even when the output weights are highly interpretable, strengthening  
454 and weakening have a very different overall behavior, and the weakening behavior is much more  
455 complex. We think that this is due to the nature of weakening: it inherently involves (an apparent)  
456 conflict between the intermediate model prediction and what the neuron promotes.

## 459 9 CONCLUSION

460 We have explored a new method for analyzing gated neurons in LLMs: computing the cosine  
461 similarities of their weights to understand their read-write functionality. Our method complements  
462 prior interpretability approaches and, though quite simple, provides striking new insights into the  
463 inner working of LLMs.

464 We have found that a large share of neurons exhibit strong RW interactions: early-middle layers are  
465 dominated by conditional strengthening neurons; weakening neurons are fewer and appear mostly in  
466 late layers. This finding is particularly noteworthy since it is universal across models, and is all the  
467 more striking since it could be observed with such a simple method.

468 Focusing on weakening neurons, a relatively small RW class, we have discovered that they have  
469 an outside impact on model behavior, including aspects as different as attribute rate (part of factual  
470 recall), and next-token entropy. We have also introduced a new analysis method, conditional ablation,  
471 which enables to find out which activations of a neuron are responsible for a given behavior. This  
472 method has shown that part of the impact of weakening neurons is due to a mechanism involving  
473 negative gate values; we are the first to observe such a mechanism.

474 Our findings open up new research questions in mechanistic interpretability, and we hope that our  
475 study will inspire further investigations. In particular, a better understanding of weakening neurons  
476 is crucial for interpreting LLMs overall. Investigating the many conditional strengthening neurons  
477 in more detail could also lead to valuable insights. In upcoming work, we plan to investigate the  
478 evolution of RW functionalities during model training. Later on, we would also like to go beyond  
479 the analysis of single neurons and address questions such as how neurons work together within and  
480 across RW classes, or whether a similar analysis also works for SAE latents.

481 <sup>12</sup>The notation is "layer.neuron", with zero-based indexing. The model – still OLMo-7B – has 32 layers, so  
482 our weakening neuron is in the final layer.

## REFERENCES

- 486  
487  
488 01.AI, :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Guoyin  
489 Wang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang  
490 Liu, Shawn Yue, Senbin Yang, Shiming Yang, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi  
491 Ren, Xinyao Niu, Pengcheng Nie, Yanpeng Li, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai,  
492 Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. Yi: Open foundation models by 01.ai, 2025. URL  
493 <https://arxiv.org/abs/2403.04652>.
- 494 Ameen Ali, Shahar Katz, Lior Wolf, and Ivan Titov. Detecting and pruning prominent but detrimental  
495 neurons in large language models. *COLM*, 2025. URL <https://arxiv.org/pdf/2507.09185>.
- 496 Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney, Stella  
497 Biderman, and Jacob Steinhardt. Eliciting latent predictions from transformers with the tuned lens,  
498 2023. URL <https://arxiv.org/pdf/2303.08112>.
- 499 Fahim Dalvi, Nadir Durrani, Hassan Sajjad, Yonatan Belinkov, Anthony Bau, and James Glass. What  
500 is one grain of sand in the desert? analyzing individual neurons in deep NLP models. *AAAI*, 2019.  
501 doi: <https://doi.org/10.1609/aaai.v33i01.33016309>.
- 502  
503 Guy Dar, Mor Geva, Ankit Gupta, and Jonathan Berant. Analyzing transformers in embedding  
504 space. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the*  
505 *61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,  
506 pp. 16124–16170, Toronto, Canada, July 2023. Association for Computational Linguistics. doi:  
507 10.18653/v1/2023.acl-long.893. URL <https://aclanthology.org/2023.acl-long.893>.
- 508  
509 Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. LLM.int8(): 8-bit matrix  
510 multiplication for Transformers at scale. *NeurIPS*, 2022. URL [https://proceedings.neurips.](https://proceedings.neurips.cc/paper_files/paper/2022/file/c3ba4962c05c49636d4c6206a97e9c8a-Paper-Confere)  
511 [nce.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/c3ba4962c05c49636d4c6206a97e9c8a-Paper-Confere).
- 512  
513 Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda  
514 Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac  
515 Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse,  
516 Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A  
517 mathematical framework for transformer circuits, 2021. URL [https://transformer-circuits.](https://transformer-circuits.pub/2021/framework/index.html)  
518 [pub/2021/framework/index.html](https://transformer-circuits.pub/2021/framework/index.html).
- 519  
520 Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec,  
521 Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish,  
522 Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposi-  
523 tion, 2022. URL [https://transformer-circuits.pub/2022/toy\\_model/index.html](https://transformer-circuits.pub/2022/toy_model/index.html).
- 524  
525 Amit Elhelo and Mor Geva. Inferring functionality of attention heads from their parameters, 2024.  
526 URL <https://arxiv.org/abs/2412.11965>.
- 527  
528 Kawin Ethayarajh. How contextual are contextualized word representations? Comparing the ge-  
529 ometry of BERT, ELMo, and GPT-2 embeddings. In Kentaro Inui, Jing Jiang, Vincent Ng, and  
530 Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Lan-  
531 guage Processing and the 9th International Joint Conference on Natural Language Processing*  
532 *(EMNLP-IJCNLP)*, pp. 55–65, Hong Kong, China, November 2019. Association for Computational  
533 Linguistics. doi: 10.18653/v1/D19-1006. URL <https://aclanthology.org/D19-1006>.
- 534  
535 Team Gemma. Gemma. 2024. doi: 10.34740/KAGGLE/M/3301. URL [https://www.kaggle.com](https://www.kaggle.com/m/3301)  
536 [/m/3301](https://www.kaggle.com/m/3301).
- 537  
538 Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are  
539 key-value memories. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-  
tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language*  
*Processing*, pp. 5484–5495, Online and Punta Cana, Dominican Republic, November 2021.  
Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.446. URL  
<https://aclanthology.org/2021.emnlp-main.446>.

- 540 Mor Geva, Avi Caciularu, Kevin Wang, and Yoav Goldberg. Transformer feed-forward layers  
541 build predictions by promoting concepts in the vocabulary space. In Yoav Goldberg, Zornitsa  
542 Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in*  
543 *Natural Language Processing*, pp. 30–45, Abu Dhabi, United Arab Emirates, December 2022.  
544 Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.3. URL  
545 <https://aclanthology.org/2022.emnlp-main.3>.
- 546 Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. Dissecting recall of factual  
547 associations in auto-regressive language models. In Houda Bouamor, Juan Pino, and Kalika Bali  
548 (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*,  
549 pp. 12216–12235, Singapore, December 2023. Association for Computational Linguistics. doi: 10  
550 .18653/v1/2023.emnlp-main.751. URL <https://aclanthology.org/2023.emnlp-main.751>.
- 551 Dirk Groeneveld, Iz Beltagy, Evan Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya  
552 Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell  
553 Authur, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel,  
554 Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal  
555 Nam, Matthew Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah,  
556 William Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan  
557 Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah Smith,  
558 and Hannaneh Hajishirzi. OLMo: Accelerating the science of language models. In Lun-Wei  
559 Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the*  
560 *Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15789–15809, Bangkok,  
561 Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-lon  
562 g.841. URL <https://aclanthology.org/2024.acl-long.841>.
- 563 Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas.  
564 Finding neurons in a haystack: Case studies with sparse probing, 2023. URL [https://arxiv.org](https://arxiv.org/pdf/2305.01610)  
565 [/pdf/2305.01610](https://arxiv.org/pdf/2305.01610).
- 566 Wes Gurnee, Theo Horsley, Zifan Carl Guo, Tara Rezaei Kheirkhah, Qinyi Sun, Will Hathaway,  
567 Neel Nanda, and Dimitris Bertsimas. Universal neurons in gpt2 language models, 2024. URL  
568 <https://arxiv.org/pdf/2401.12181>.
- 569 Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with gaussian  
570 error linear units. *CoRR*, 2016. URL <http://arxiv.org/abs/1606.08415>.
- 571 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
572 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,  
573 L lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas  
574 Wang, Timoth e Lacroix, and William El Sayed. Mistral 7b, 2023. URL [https://arxiv.org/ab](https://arxiv.org/abs/2310.06825)  
575 [s/2310.06825](https://arxiv.org/abs/2310.06825).
- 576 Subhash Kantamneni, Joshua Engels, Senthoooran Rajamanoharan, Max Tegmark, and Neel Nanda.  
577 Are sparse autoencoders useful? a case study in sparse probing. *arXiv*, 2025. URL [https:](https://arxiv.org/pdf/2502.16681.pdf)  
578 [/arxiv.org/pdf/2502.16681.pdf](https://arxiv.org/pdf/2502.16681.pdf).
- 579 Linghao Kong, Angelina Ning, Micah Adler, and Nir Shavit. Negative pre-activations differentiate  
580 syntax. *arXiv*, 2025. URL <https://arxiv.org/pdf/2509.24198.pdf>.
- 581 Olga Kovaleva, Saurabh Kulshreshtha, Anna Rogers, and Anna Rumshisky. BERT busters: Outlier  
582 dimensions that disrupt transformers. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto  
583 Navigli (eds.), *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp.  
584 3392–3405, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/  
585 2021.findings-acl.300. URL <https://aclanthology.org/2021.findings-acl.300>.
- 586 Vedang Lad, Wes Gurnee, and Max Tegmark. The remarkable robustness of llms: Stages of inference?,  
587 2024. URL <https://arxiv.org/abs/2406.19384>.
- 588 Patrick Leask, Bart Bussmann, Michael Pearce, Joseph Bloom, Curt Tigges, Noura Al Moubayed,  
589 Lee Sharkey, and Neel Nanda. Sparse autoencoders do not find canonical units of analysis. *arXiv*,  
590 2025. URL <https://arxiv.org/pdf/2502.04878.pdf>.
- 591
- 592
- 593

- 594 Minhyeok Lee. Mathematical analysis and performance evaluation of the GELU activation function  
595 in deep learning. *Journal of Mathematics*, 2023. doi: <https://doi.org/10.1155/2023/4229924>.  
596
- 597 Joseph Miller and Clement Neo. We found an neuron in gpt-2, 2023. URL <https://www.lesswrong.com/posts/cgqh99SHsCv3jJYDS/we-found-an-neuron-in-gpt-2>.  
598
- 599 Beren Millidge and Sid Black. The singular value decompositions of transformer weight matrices are  
600 highly interpretable, 2022. URL [https://www.lesswrong.com/posts/mkbGjzx8d8XqKHzA/](https://www.lesswrong.com/posts/mkbGjzx8d8XqKHzA/the-singular-value-decompositions-of-transformer-weight)  
601 [the-singular-value-decompositions-of-transformer-weight](https://www.lesswrong.com/posts/mkbGjzx8d8XqKHzA/the-singular-value-decompositions-of-transformer-weight).  
602
- 603 Ari S. Morcos, David G.T. Barrett, Neil C. Rabinowitz, and Matthew Botvinick. On the importance  
604 of single directions for generalization, 2018. URL <https://arxiv.org/pdf/1803.06959.pdf>.
- 605 Aaron Mueller, Atticus Geiger, Sarah Wiegrefe, Dana Arad, Iván Arcuschin, Adam Belfki, Yik Siu  
606 Chan, Jaden Fiotto-Kaufman, Tal Haklay, Michael Hanna, Jing Huang, Rohan Gupta, Yaniv  
607 Nikankin, Hadas Orgad, Nikhil Prakash, Anja Reusch, Aruna Sankaranarayanan, Shun Shao,  
608 Alessandro Stolfo, Martin Tutek, Amir Zur, David Bau, and Yonatan Belinkov. MIB: a mechanistic  
609 interpretability benchmark. *arXiv*, 2025. URL <https://arxiv.org/pdf/2504.13151>.
- 610 Neel Nanda. Neuroscope. Website, 2022. URL <https://neuroscope.io>.  
611
- 612 Neel Nanda and Joseph Bloom. Transformerlens. [https://github.com/TransformerLensOrg/Tr](https://github.com/TransformerLensOrg/TransformerLens)  
613 [ansformerLens](https://github.com/TransformerLensOrg/TransformerLens), 2022.  
614
- 615 Jingcheng Niu, Andrew Liu, Zining Zu, and Gerald Penn. What does the knowledge neuron thesis  
616 have to do with knowledge?, 2024. URL <https://arxiv.org/pdf/2405.02421>.
- 617 nostalgebraist. Interpreting gpt: The logit lens, 2020. URL [https://www.lesswrong.com/posts/](https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens)  
618 [AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens](https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens).  
619
- 620 Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry  
621 of large language models. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller,  
622 Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International*  
623 *Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp.  
624 39643–39666. PMLR, 21–27 Jul 2024. URL [https://proceedings.mlr.press/v235/park2](https://proceedings.mlr.press/v235/park24c.html)  
625 [4c.html](https://proceedings.mlr.press/v235/park24c.html).
- 626 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language  
627 models are unsupervised multitask learners. 2019.
- 628 Prajit Ramachandran, Barret Zoph, and Quoc V. Le. Searching for activation functions. *arXiv*, 2017.  
629 URL <https://arxiv.org/pdf/1710.05941>.  
630
- 631 Cody Rushing and Neel Nanda. Explorations of self-repair in language models. In Ruslan Salakhut-  
632 dinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix  
633 Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, vol-  
634 ume 235 of *Proceedings of Machine Learning Research*, pp. 42836–42855. PMLR, 21–27 Jul  
635 2024. URL <https://proceedings.mlr.press/v235/rushing24a.html>.
- 636 Lee Sharkey, Dan Braun, and Beren Millidge. [interim research report] taking features out of  
637 superposition with sparse autoencoders, 2022. URL [https://www.alignmentforum.org/posts](https://www.alignmentforum.org/posts/z6QQJbtpkEAX3Aojj/interim-research-report-taking-features-out-of-superposition)  
638 [/z6QQJbtpkEAX3Aojj/interim-research-report-taking-features-out-of-superposi](https://www.alignmentforum.org/posts/z6QQJbtpkEAX3Aojj/interim-research-report-taking-features-out-of-superposition)  
639 [tion](https://www.alignmentforum.org/posts/z6QQJbtpkEAX3Aojj/interim-research-report-taking-features-out-of-superposition).  
640
- 641 Noam Shazeer. Glu variants improve transformer, 2020. URL [https://arxiv.org/pdf/2002.052](https://arxiv.org/pdf/2002.05202)  
642 [02](https://arxiv.org/pdf/2002.05202).
- 643 Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur,  
644 Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, Ananya Jha,  
645 Sachin Kumar, Li Lucy, Xinxin Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas  
646 Muennighoff, Aakanksha Naik, Crystal Nam, Matthew Peters, Abhilasha Ravichander, Kyle  
647 Richardson, Zejiang Shen, Emma Strubell, Nishant Subramani, Oyvind Tafjord, Evan Walsh,  
Luke Zettlemoyer, Noah Smith, Hannaneh Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge,

- 648 and Kyle Lo. Dolma: an open corpus of three trillion tokens for language model pretraining  
649 research. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd*  
650 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp.  
651 15725–15788, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi:  
652 10.18653/v1/2024.acl-long.840. URL <https://aclanthology.org/2024.acl-long.840>.
- 653 Alessandro Stolfo, Ben Wu, Wes Gurnee, Yonatan Belinkov, Xingyi Song, Mrinmaya Sachan, and  
654 Neel Nanda. Confidence regulation neurons in language models, 2024. URL <https://arxiv.org/abs/2406.16254>.
- 655  
656  
657 Mingjie Sun, Xinlei Chen, J. Zico Kolter, and Zhuang Liu. Massive activations in large language  
658 models. *arXiv*, 2024. URL <https://arxiv.org/pdf/2402.17762>.
- 659 William Timkey and Marten van Schijndel. All bark and no bite: Rogue dimensions in transformer  
660 language models obscure representational quality. In Marie-Francine Moens, Xuanjing Huang,  
661 Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical*  
662 *Methods in Natural Language Processing*, pp. 4527–4546, Online and Punta Cana, Dominican  
663 Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.e  
664 mnlp-main.372. URL <https://aclanthology.org/2021.emnlp-main.372>.
- 665 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
666 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand  
667 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language  
668 models. *ArXiv*, abs/2302.13971, 2023. URL <https://arxiv.org/pdf/2302.13971>.
- 669 Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
670 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Neural Information Processing*  
671 *Systems*, 2017. URL <https://api.semanticscholar.org/CorpusID:13756489>.
- 672 Roman Vershynin. *High dimensional probability*. 2nd edition, 2025. URL [https://www.math.uci](https://www.math.uci.edu/~rvershyn/papers/HDP-book/HDP-2.pdf)  
673 [.edu/~rvershyn/papers/HDP-book/HDP-2.pdf](https://www.math.uci.edu/~rvershyn/papers/HDP-book/HDP-2.pdf).
- 674 Elena Voita, Javier Ferrando, and Christoforos Nalmpantis. Neurons in large language models: Dead,  
675 n-gram, positional. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the*  
676 *Association for Computational Linguistics: ACL 2024*, pp. 1288–1301, Bangkok, Thailand, August  
677 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.75. URL  
678 <https://aclanthology.org/2024.findings-acl.75>.
- 679 Zhengxuan Wu, Aryaman Arora, Atticus Geiger, Zheng Wang, Jing Huang, Dan Jurafsky, Christo-  
680 pher D. Manning, and Christopher Potts. AxBench: Steering LLMs? even simple baselines  
681 outperform sparse autoencoders. *ICML*, 2025. URL <https://arxiv.org/abs/2501.17148>.
- 682 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li,  
683 Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang,  
684 Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai,  
685 Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng  
686 Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai  
687 Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan  
688 Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang  
689 Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2  
690 technical report. *arXiv preprint arXiv:2407.10671*, 2024.
- 691 Xiutian Zhao, Rochelle Choenni, Rohit Saxena, and Ivan Titov. Finding culture-sensitive neurons in  
692 vision-language models. *arXiv*, 2025. URL <https://arxiv.org/pdf/2510.24942>.
- 693

## 697 A LIMITATIONS

698 We focus on a *parameter-based* interpretation of *single neurons*. This has the advantage of being  
699 simple and efficient, but is also inherently limited in scope. Accordingly, our method is not designed  
700 to replace other neuron analysis methods, but to complement them.  
701

The mathematical similarities of weights are insightful, but they should not be taken as one-to-one representations of semantic similarity. We find cases in which close-to-orthogonal vectors represent very similar concepts (double checking, section H).

Finally, while our findings are highly significant and relevant, we do not yet fully understand the reasons behind them.

## B DEONTOLOGY STATEMENTS

### B.1 LICENSES AND LANGUAGES OF MODELS AND DATA

**Gemma.** To download the model one needs to explicitly accept the terms of use. NLP research is explicitly listed as an intended usage. Primarily English and code Gemma (2024).

**Llama.** Inference code and weights under an ad hoc license. There is also an “Acceptable Use Policy”. Our work is well within those terms. Languages mostly include English and programming languages, but also Wikipedia dumps from “bg, ca, cs, da, de, en, es, fr, hr, hu, it, nl, pl, pt, ro, ru, sl, sr, sv, uk” Touvron et al. (2023).

**OLMo and Dolma.** Training and inference code, weights (OLMo), and data (Dolma) under Apache 2.0 license. “The Science of Language Models” is explicitly mentioned as an intended use case. Dolma is quality-filtered and designed to contain only English and programming languages (though we came across some French sentences as well, see table 3) Groeneveld et al. (2024); Soldaini et al. (2024).

**Mistral.** Inference code and weights are released under the Apache 2.0 license, but accessing them requires accepting the terms. Languages are not explicitly mentioned in the paper, but clearly include English and code Jiang et al. (2023).

**Qwen.** Inference code and weights under Apache 2.0 license. Supports “over 29 languages, including Chinese, English, French, Spanish, Portuguese, German, Italian, Russian, Japanese, Korean, Vietnamese, Thai, Arabic, and more” Yang et al. (2024).

**Yi.** Inference code and weights under Apache 2.0 license. Trained on English and Chinese 01.AI et al. (2025).

### B.2 COMPUTATIONAL COMPLEXITY

All our experiments can be run on a single NVIDIA RTX A6000 (48GB). We use TransformerLens Nanda & Bloom (2022). A colleague kindly provided us with a version that also supports OLMo.

The main analysis, computing the weight cosines, needs less than a minute per model.

Other parts were more expensive:

- For the ablations (section 6), each run on Dolma used took approximately 8 hours.
- For the activation-based analysis in section 8, we needed a single run of  $\approx 25$  h to store the max/min activating examples for all neurons, and then  $\approx 45$  s per neuron ( $\approx 5$  min) to recompute its activations on the relevant texts and visualize them.
- Another expensive part is computing the randomness regions based on mismatched cosines (sections 4.3 and E). The time complexity is  $O(n^2)$  in the number of neurons per layer, since we have to consider every pair of neurons. Since however we found that this baseline is hardly different from the more “naive” Gaussian one, we suggest that future work could just leave out this step.
- Finally, our weight processing makes model loading last about a minute. A possible solution in future work would be to save the preprocessed weights.

### B.3 LLM USE

We used LLM assistants to help with programming.

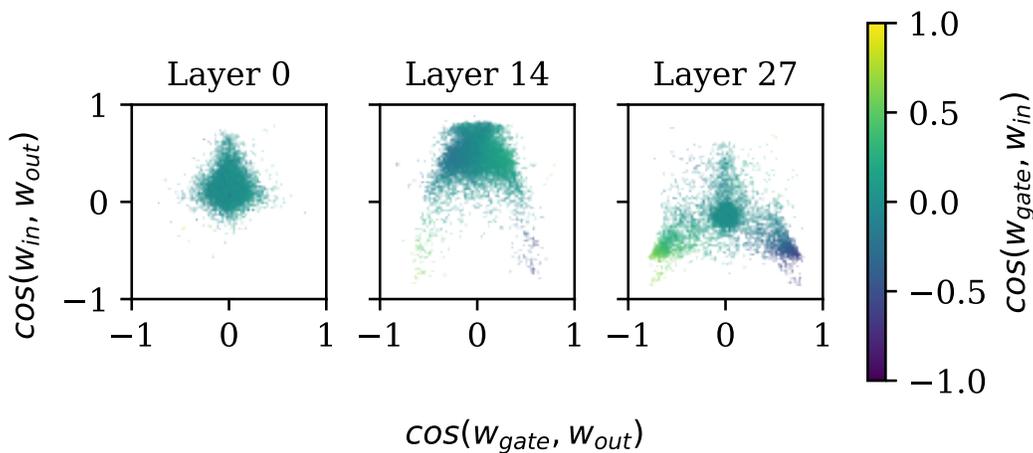


Figure 5: Equivalent of figure 52, but *without* weight processing. (We don't include the randomness regions.)

## C WEIGHT PREPROCESSING

In this section we argue for our weight preprocessing step described in section 3.2: multiplying  $w_{in}$  and  $w_{out}$  by the sign of  $\cos(w_{gate}, w_{in})$ .

First of all, this processing **does not affect model behavior**: In equation (2), if we replace  $w_{in}$  by  $-w_{in}$  and  $w_{out}$  by  $-w_{out}$ , the two minus signs cancel out. We call this the **symmetry property** of gated activation functions.

We find neurons slightly **easier to interpret** after applying this weight preprocessing step, for the following reasons:

- The two reading weight vectors,  $w_{gate}$  and  $w_{in}$ , now always have a non-negative cosine similarity, i.e., they do not point in opposite directions. This makes it easier to reason about what causes a neuron to activate (there are less minus signs to worry about). This is especially relevant for the case studies (section 8).
- In scatter plots like figure 2, neurons that belong together are in the same area of the plot. Consider any dot in figure 2 (which was made *with* the weight processing), for example on the bottom left of the subplot (i.e. a weakening neuron). If we undo the weight processing,  $\cos(w_{gate}, w_{out})$  (say  $-0.8$ ) may be randomly replaced by its opposite (say  $+0.8$ ), in which case the same neuron would now appear on the bottom *right*. Accordingly, in the equivalent figure *without* weight processing (figure 5), we see two distinct clusters of weakening neurons, at the bottom left and bottom right, while these clusters are not different in terms of neuron behavior.
- Similarly, for the conditional ablations introduced in section 6.2, the four cases correspond to real distinctions. Without weight preprocessing, equivalent cases would be more complicated to define (e.g. "gate+\_post+" would be defined as " $x_{gate} > 0$  and  $x_{post}$  has the same sign as  $\cos(w_{gate}, w_{in})$ ").

## D FROM NEURONS TO LINEAR COMBINATIONS

In section 4.1 we claimed that neuron-based findings can carry over to linear combinations of neurons, at least to some extent.

We will give the argument for one example: a combination of two prototypical strengthening neurons (let's call them  $n^{(1)}$  and  $n^{(2)}$ ). Let their weights be  $w_{gate}^{(1)} = w_{in}^{(1)} = w_{out}^{(1)} := w^{(1)}$  and

810  $\mathbf{w}_{\text{gate}}^{(2)} = \mathbf{w}_{\text{in}}^{(2)} = \mathbf{w}_{\text{out}}^{(2)} := \mathbf{w}^{(2)}$ . We consider the multi-neuron feature  $n^{(1)} + n^{(2)}$  (defined by  $n^{(1)}$   
811 and  $n^{(2)}$  activating together with the same strengths).  
812

813 This feature will activate when  $\mathbf{x}_{\text{norm}}$  contains both  $\mathbf{w}^{(1)}$  and  $\mathbf{w}^{(2)}$ , which implies containing  $\mathbf{w}^{(1)} +$   
814  $\mathbf{w}^{(2)}$  (though the converse does not hold). When activating, it will write a positive multiple to  
815  $\mathbf{w}^{(1)} + \mathbf{w}^{(2)}$  to the residual stream. Thus this linear combination of strengthening neurons has a  
816 strengthening-like behavior itself.

817 Similar arguments can be made for any pair of neurons of the same RW class, and also for more than  
818 two neurons.  
819

820 In particular, in a layer containing many conditional strengthening neurons (like early-middle layers  
821 of all models), a sparse linear combination of neurons is reasonably likely to consist of conditional  
822 strengthening neurons, in which case it will display a behavior similar to conditional strengthening.  
823

## 824 E RANDOM BASELINES

825 Here we describe our two baselines: random initialization and mismatched cosines.  
826

827 In a randomly initialized model, all cosine similarities would be very close to zero: In  $n$  dimensions,  
828 absolute cosine similarities behave like  $1/\sqrt{n}$  (Vershynin, 2025, p. 68). More precisely, the cosines  
829 follow a beta distribution with parameters  $(d_{\text{model}} - 1)/2, (d_{\text{model}} - 1)/2$ , rescaled to the range  
830  $[-1, 1]$ .<sup>13</sup> Taking, e.g.,  $d_{\text{model}} = 4096$  (as e.g. in OLMo-7B), we get a 95% randomness range  
831 of approximately  $[-0.03, 0.03]$ . This is empirically confirmed on the first training checkpoint of  
832 OLMo-7B-0424 (figure 6).

833 Inspired by work on outlier dimensions in the *activations* of Transformers (Ethayarajh, 2019; Kovaleva  
834 et al., 2021; Timkey & van Schijndel, 2021; Dettmers et al., 2022; Sun et al., 2024), we suspected  
835 that a similar phenomenon might be at work in the *weights*, making cosine similarities artificially  
836 high. To account for this possibility, we construct a second baseline specific to each model layer: We  
837 compute all the (e.g.)  $\cos(\mathbf{w}_{\text{in}}, \mathbf{w}_{\text{out}})$  of a layer, even if the two weights belong to different neurons.  
838 If a cosine similarity is higher than most of these mismatched cosines, it is likely not due to an outlier  
839 dimension common to all neurons of the layer, but reflects something specific to this neuron.  
840

## 841 F ABLATION EXPERIMENTS

### 842 F.1 HYPOTHESES AND CHOICE OF METRICS

843 We originally had two hypotheses (which turned out to be wrong, see section 6):  
844

- 845 • We hypothesized that *conditional strengthening* neurons might contribute to *subject en-*  
846 *richment* Geva et al. (2023), a crucial step of factual recall that involves MLPs writing  
847 appropriate attributes for the given subject. Both phenomena occur in roughly the same  
848 layers, and similar  $\mathbf{w}_{\text{in}}$  and  $\mathbf{w}_{\text{out}}$  could correspond to related concepts.  
849
- 850 • We expected that *weakening* neurons would make the output distribution flatter, i.e. *increase*  
851 *the entropy*. This could happen by reducing the probability of high-ranking tokens (weaken-

852 ing directions corresponding to tokens) or by increasing the probability of very low-ranking  
853 tokens (weakening directions corresponding to negations of tokens).  
854

855 This is why we tested the two metrics of *attribute rate* (a proxy of subject enrichment) and *entropy*.  
856 We additionally considered the *loss*, and, following Gurnee et al., 2024 analysis of entropy neurons,  
857 *rank* of the correct token and *scale* of the final hidden state.  
858

### 859 F.2 NUMBER OF NEURONS TO ABLATE

860 In preliminary experiments, we tried ablating 24 or 243 neurons in each run, which is the number of  
861 strengthening or weakening neurons in OLMo, respectively. Ablating 24 neurons (of any class) did  
862

863 <sup>13</sup><https://stats.stackexchange.com/questions/85916/distribution-of-scalar-products-of-two-random-unit-vectors-in-d-dimensions>

864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

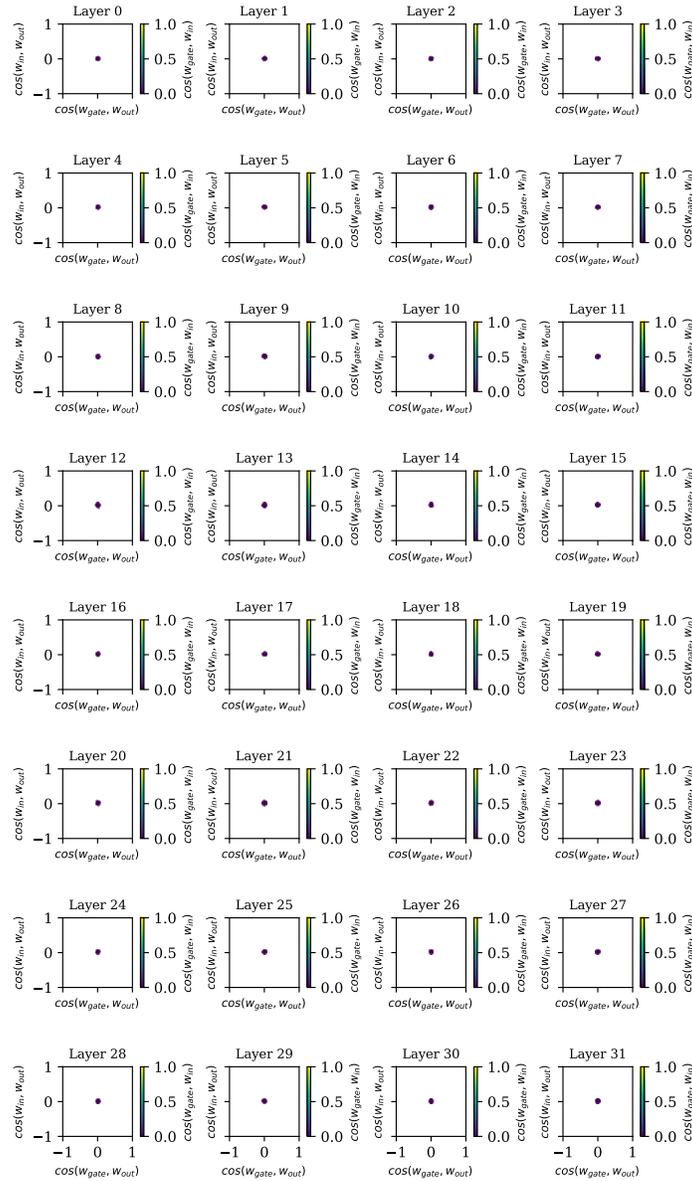


Figure 6: Equivalent of figure 2 for the randomly initialized OLMo-7B model (training checkpoint 0). Whatever doesn't look like this, is significant.

not have a clear impact on any metric (including loss), so we stuck with 243, which has a visible impact on the loss, but does not break the model altogether. For this reason we exclude strengthening neurons from most of our experiments, since there are only 24 of them.

### F.3 DETAILS ON ATTRIBUTE RATE

Our investigation of attribute rate closely follows Geva et al. (2023). It requires a dataset of subject-attribute mappings that we didn't have access to. In order to replicate this dataset, we closely followed the procedure described in their paper, which assumes attributes are tokens that appear in the same Wikipedia paragraph as the subject (excluding stopwords). We used the Wikipedia dump from October 20, 2021, instead of October 13, since there is an official dump made at this date.<sup>14</sup> To improve replicability, we publish our complete code as well as our subject-attribute dataset.

### F.4 MEAN ABLATION

#### F.4.1 METHOD

**Computing the means.** We pre-computed the mean activation of every neuron on the same 20M token subset of Dolma that we also used for the actual ablation experiment.

**Conditional ablations.** For conditional ablations, we replace the neuron activation by the mean value it would have in the corresponding case (not the mean activation of the neuron overall). For example, in the case "gate-\_post+":

- we replace the activation ( $x_{\text{post}}$ ) only when the condition "gate-\_post+" is fulfilled, i.e. when  $x_{\text{gate}} < 0$  and  $x_{\text{post}} > 0$  – this is just the definition of conditional ablation;
- the value that we replace it with is the mean value of  $x_{\text{post}}$  across the cases in which  $x_{\text{gate}} < 0$  and  $x_{\text{post}} > 0$ .

#### F.4.2 RESULTS

We show the results in figures 10 to 13, 25 to 32 and 36 to 39.

Regarding **attributes rate**, we find again that non-weakening neurons have no effect (figures 11 to 13 and 25). Mean-ablating weakening neurons has a somewhat inconsistent effect (figure 10): In middle layers the intervention reduces attributes rate (though less starkly than with zero ablation), suggesting that these neurons play a role in subject enrichment. In some late layers however (especially layer 29), the intervention *increases* attributes rate, suggesting that the neurons help avoiding an explosion of the attributes rate at this computation stage. (This would make sense: at the end of the forward pass the model wants to predict the next token, not some attribute of the subject.)

Regarding **entropy**, it remains the case that ablating weakening neurons has the biggest effect (figure 29). (In both zero and ablation, strengthening neurons also often reduce entropy, but this is expected. See figures 17 and 25) In particular, there is still a substantial number of cases in which weakening neurons *reduce* entropy by about 10 nats. However, entropy changes are now more evenly distributed: there is also a substantial (only slightly smaller) amount of cases in which weakening neurons *increase* entropy, which aligns better with their expected behavior. **Conditional ablations** (figure 36) partially recover both effects with the case "gate-\_post+" (as in zero ablation), but also with the case "gate+\_post-" (contrary to zero ablation).

Mean ablation thus recovers effects of weakening neurons that would go unnoticed when just using zero ablation. We hypothesize that these additional effects happen when activations are relatively close to zero but far away from their mean. This remains to be tested in future work.

## G ACTIVATION FREQUENCIES

This section extends section 7.

Complete results are in table 5 and figures 7 to 9.

<sup>14</sup>A list of dumps by date is available at <https://archive.org/search?query=subject>.

Table 2: Overview of prediction/suppression neurons chosen for case studies in section 8

Neuron	RW category	$\cos(\mathbf{w}_{\text{gate}}, \mathbf{w}_{\text{in}})$	$\cos(\mathbf{w}_{\text{gate}}, \mathbf{w}_{\text{out}})$	$\cos(\mathbf{w}_{\text{in}}, \mathbf{w}_{\text{out}})$
28.4737	strengthening	0.5290	0.5048	0.7060
28.9766	conditional strengthening	0.4764	0.4119	0.5982
31.9634	weakening	-0.7164	0.7218	-0.8542
29.10900	conditional weakening	0.4988	-0.4992	-0.5775
30.10972	proportional change	-0.4543	0.5814	-0.4182
29.4180	orthogonal output	-0.0272	-0.4057	0.0669

The last layer displays a different pattern than the rest (last subplot in figure 7). Here the correlation is positive (+0.29), and we can distinguish two clusters of neurons: One cluster has a medium-negative  $\cos(\mathbf{w}_{\text{in}}, \mathbf{w}_{\text{out}})$  (around  $-0.3$ ) and activates very rarely; another one is much more spread out (both in terms of  $\cos(\mathbf{w}_{\text{in}}, \mathbf{w}_{\text{out}})$  and activation frequency), centers at a weaker negative cosine similarity ( $-0.1$  to  $-0.2$ ) and activates a bit more than half of the time. The presence of these two clusters leads to the slightly positive correlation. Comparing with the other plots suggests that the first cluster mostly corresponds to weakening neurons and atypical proportional change neurons.

We do not find such striking patterns with gate-out or gate-in similarities.

## H DOUBLE CHECKING

In our case studies, we observe that many neurons have the property of **double checking**: The two reading weight vectors ( $\mathbf{w}_{\text{gate}}$  and  $\mathbf{w}_{\text{in}}$ ) are approximately orthogonal, but still intuitively represent the same concept.

We characterize double checking as follows: The sets of meaningful vectors *most similar* to  $\mathbf{w}_{\text{gate}}$  and  $\mathbf{w}_{\text{in}}$  have a high overlap. More formally, let  $U = \{u_0, \dots, u_{d_{\text{vocab}}}\}$  be the set of unembedding vectors; then

$$\arg \max_{u \in U} \cos(u, \mathbf{w}_{\text{gate}}) \approx \arg \max_{u \in U} \cos(u, \mathbf{w}_{\text{in}}).$$

This phenomenon is *possible* because random vectors in high dimensions are “lone stars” (Vershynin, 2025, p. 68). If this is the case for the unembedding vectors, it is plausible that we can find  $\mathbf{w}_{\text{gate}}, \mathbf{w}_{\text{in}}$  that are reasonably similar to a  $u_i$  but not to any other  $u_j$ . These  $\mathbf{w}_{\text{gate}}, \mathbf{w}_{\text{in}}$  can even be (approximately) orthogonal to each other, as in the following three-dimensional toy example:  $\mathbf{w}_1 = (1, 0, 0)$ ,  $\mathbf{w}_2 = (0, 1, 0)$ ,  $\mathbf{w}_{\text{gate}} = (1, 0, 1)$ ,  $\mathbf{w}_{\text{in}} = (1, 0, -1)$ .

However the phenomenon is *unlikely* to occur in random vectors, and hence is a significant finding: If choosing  $\mathbf{w}_{\text{gate}}, \mathbf{w}_{\text{in}}$  randomly, we would expect them to be approximately orthogonal to *all* unembedding vectors; and even if both were somewhat similar to an unembedding, we certainly wouldn’t expect it to be the same unembedding for both.

We would also not naively expect this phenomenon in a trained network: If the role of both  $\mathbf{w}_{\text{gate}}$  and  $\mathbf{w}_{\text{in}}$  is to detect a concept (e.g. a token prediction) represented by a vector  $u$ , then we would get the best performance with  $\mathbf{w}_{\text{gate}} = \mathbf{w}_{\text{in}} = u$ , i.e.,  $\mathbf{w}_{\text{gate}}, \mathbf{w}_{\text{in}}$  would not be orthogonal.

Double checking is therefore likely to be a useful feature for the model. We hypothesize that this is because it shrinks the region in model space that activates the neuron positively. If (say)  $\mathbf{w}_{\text{in}} = \mathbf{w}_{\text{gate}} = (1, 0)$ , the neuron activates whenever the (normalized) residual input  $x$  satisfies  $x \cdot (1, 0) > 0$ ; this happens on the whole half-space  $x_1 > 0$ . If however  $\mathbf{w}_{\text{gate}} = (1, 0)$  and  $\mathbf{w}_{\text{in}} = (0, 1)$ , the neuron activates positively only in the first quadrant ( $x_1, x_2 > 0$ ).

This behavior thus enables more precise concept detection. This may explain why conditional neurons are more frequent than their unconditional counterparts.

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

Table 3: Description of the weight vectors of the selected *prediction* neurons, by top tokens or similarity to  $w_{out}$ . The question mark, ?, signals unknown unicode characters. The last column presents the (shortened) text samples on which the respective neuron activates most strongly (positively or negatively).

Neuron, RW class	$w_{gate}$		$w_{in}$		$w_{out}$	Top activations
28.4737 strengthening	$\approx w_{out}$		$\approx w_{out}$		pos: <i>review</i> neg: <i>view</i>	pos (13.75): <i>Download EBOOK [...] Description of the book [...] \n -&gt; Reviews</i> neg (-2.25): <i>The answer's at the bottom of <b>this</b> -&gt; post</i>
28.9766 conditional strengthening	pos: <i>well</i> neg: <i>far</i> <i>high</i>		$\approx w_{out}$		pos: <i>well</i> neg: <i>well</i>	pos (18.63): <i>Could have saved myself some time. <b>Oh</b> -&gt; , well</i> neg (-3.66): <i>Seek to understand them <b>more</b> -&gt; fully</i>
31.9634 weakening	$\approx -w_{out}$		$\approx -w_{out}$		pos: <i>again</i> neg: <i>Again</i>	pos (3.48): <i>the areas of the doorjamb where the <b>door</b> -&gt; often</i> neg (-5.12): <i>jumping off the roof of his Los Angeles apartment building. -&gt; <b>Meanwhile</b></i>
29.10900 conditional weakening	pos: <i>today</i> neg: <i>these</i> <i>these</i> <i>nowa-</i> <i>days</i>		$\approx -w_{out}$		pos: <i>these</i> neg: <i>These</i>	pos (12.79): <i>social media tools change and come and go at the drop of a <b>hat</b> -&gt; .</i> neg (-2.18): <i>la couleur de sa robe <b>et</b> -&gt; le</i>
30.10972 proportional change	$\approx w_{out}$		pos: <i>when</i> neg: <i>timing</i> <i>dates</i>		neg: <i>when</i> neg: <i>when</i>	pos (6.14): <i>puts you on multiple webpages <b>at</b> -&gt; as soon as</i> neg (-2.67): <i>Take pleasure in the rest of the new year. -&gt; <b>You</b></i>
29.4180 orthogonal output	pos: <i>here</i> neg: <i>there</i> <i>therein</i> <i>we</i>		pos: ?	neg: <i>here</i> <i>in</i>	neg: <i>there</i> neg: <i>there</i>	pos (2.31): <i>without any consideration being issued or paid <b>there</b> -&gt; for</i> neg (-14.41): <i>here <b>or</b> -&gt; there</i>

Table 4: Description of the weight vectors of the selected *prototypical* neurons, by top tokens or similarity to  $w_{out}$ . The question mark, ?, signals unknown unicode characters. The last column presents the (shortened) text samples on which the respective neuron activates most strongly (positively or negatively).

Neuron, RW class	$w_{gate}$	$w_{in}$	$w_{out}$	Top activations
25.9997 strengthening	$\approx w_{out}$	$\approx w_{out}$	pos: <i>S</i> neg: <i>S</i> <i>Choco-late</i> <i>Cour</i>	pos (10.45): <i>when Stannis gives his opinion on Janos -&gt; Slynt</i> neg (-0.72): <i>said owner Hieu Than, -&gt; who</i>
5.10602 conditional strengthening	pos: <i>t</i> neg: <i>deep</i> <i>as</i> <i>hum</i>	$\approx w_{out}$	pos: <i>as</i> neg: <i>ating</i> <i>t</i> <i>their</i>	pos (1.56): <i>a big-time rocker, playing aren -&gt; as</i> neg (-0.68): <i>workers don -&gt; 't</i>
31.7117 weakening	$\approx -w_{out}$	$\approx -w_{out}$	pos: <i>by</i> neg: <i>ani</i> <i>by</i> <i>iw</i>	pos (7.27): <i>deux projets de décision figurant -&gt; dans</i> neg (-5.23): <i>Take pleasure in the rest of the new year. -&gt; You</i>
23.6543 conditional weakening	pos: <i>the</i> neg: <i>ham</i> <i>a</i> <i>aden</i>	$\approx -w_{out}$	pos: <i>Op</i> neg: <i>rom</i> <i>AB</i> <i>c</i>	pos (1.06): <i>High-value and sub-high-value -&gt; areas</i> neg (-0.99): <i>:basis and -&gt; :sub-category</i>
25.7415 proportional change	$\approx w_{out}$	pos: <i>berry</i> neg: <i>a</i> <i>rod</i> <i>the</i>	pos: ? neg: <i>Nine</i> <i>Hart</i> <i>jin</i>	pos (2.36): <i>180 °C ( -&gt; 350 °F)</i> neg (-1.55): <i>// @Component({\n // -&gt; selector</i>

## I CASE STUDIES

### I.1 NEURON CHOICE

We used two different methods to find interesting neurons:

**First**, we selected among *prediction neurons* in the sense of Gurnee et al. (2024). These are defined as neurons whose  $\cos(W_U, w_{out})$  has a high kurtosis; in other words, they boost predictions of a small set of tokens while leaving other token scores virtually unchanged. Specifically, from each discrete RW class we chose the neuron with the highest kurtosis. This first method guarantees finding interpretable neurons in terms of output behavior, though not necessarily an interpretable *overall* behavior. See table 2 for an overview of neurons chosen by this method.

A downside is that prediction neurons tend to appear in later layers only. Therefore this neuron choice does not help understand what happens in early layers, especially why there are so many conditional strengthening neurons. We therefore also use a **second** method: We just select the most prototypical neuron from each class. For example, for conditional strengthening, we take the neuron with the highest  $\cos(w_{in}, w_{out})$  among those neurons whose  $\cos(w_{gate}, w_{out})$  is within the randomness range (sections 4.3 and E). This method led to choosing the neurons 5.10602 (conditional strengthening), 23.6543 (conditional weakening), 25.7415 (proportional change), 25.9997 (strengthening), 31.7117 (weakening).

### I.2 METHODS

Additionally to our RW analysis, we use two well-established neuron analysis methods:

First, we project neuron weights to vocabulary space with the unembedding matrix  $W_U$  and inspect high-scoring tokens.

Second, we find text examples on which the neurons are strongly activated (positively or negatively). For each neuron we save the 16 strongest positive and negative activations, respectively.

### I.3 DETAILED ANALYSIS OF WEAKENING NEURON 31.9634

Here we say a bit more about the neuron analyzed in section 8.

Judging by the weights, we would predict the following: The neuron activates positively when the residual stream contains the “minus *again*” direction, and then weakens that direction by writing “plus *again*”. The neuron activates negatively when the residual stream contains information both for and against predicting *again*, and then weakens the *again* direction. Given that  $w_{\text{gate}}$  and  $w_{\text{in}}$  are highly similar ( $\cos(w_{\text{gate}}, w_{\text{in}}) = 0.7164$ ), we would expect that it is easier for the neuron to activate positively (with  $x_{\text{gate}}$  and  $x_{\text{in}}$  of the same sign).

When actually recording activations of the neuron, we get a more complex picture: First of all, the neuron often activates negatively. Strong negative activations are often on punctuation, and the actual next token is often *meanwhile* or *instead* (and not *again*). On the positive side, the strongest activations do not have any obvious semantic relationship to *again*. We also observed weaker positive activations when *again* is a plausible continuation, e.g., on the token *once* (as in *once again*). These are cases with negative  $x_{\text{gate}}$  values (and also  $x_{\text{in}} < 0$ , hence positive activations) – a case that we found to be important in section 6.2. In these cases, *again* is already weakly present in the residual stream before the last MLP, and the neuron reinforces *again*.

Thus the behavior of this particular weakening neuron is interpretable in the  $x_{\text{gate}} < 0$  case, echoing our finding from section 6.2 that this case is surprisingly relevant to model behavior. The  $x_{\text{gate}} > 0$  case is less interpretable for this particular neuron, even though this case is more frequent and can lead to stronger activations. Nevertheless we have some hypotheses for the strong activations as well: For strong positive activations (which showed no clear pattern), we hypothesize that sometimes the residual stream ends up near “minus *again*” for semantically unrelated reasons (there are many more possible concepts than dimensions, so the corresponding directions cannot be fully orthogonal; see Elhage et al., 2022); in these cases the neuron would reduce the unjustified presence of this “minus *again*” direction. With strong negative activations (where the next token was often *meanwhile* or *instead*), the neuron may ensure only these tokens are predicted, and not the relatively similar *again*.

### I.4 RESULTS AND ANALYSIS FOR PREDICTION NEURONS

See table 3.

**Strengthening neuron 28.4737**<sup>15</sup> predicts *review* (and related tokens) if activated positively, which happens if *review* is already present in the residual stream. The maximally positive activations are in standard contexts that continue with *review* or similar, such as the newline after the description of an e-book (the next paragraph often is the beginning of a review).

**Conditional strengthening neuron 28.9766**’s RW functionality concerns *well* and similar tokens. 28.9766 promotes them if activated positively, which happens when both  $w_{\text{gate}}$  and  $w_{\text{in}}$  indicate that *well* is represented in the residual stream. This is a case of double checking. The maximally positive activation in our sample occurs on *Oh*, in a context in which *Oh, well* makes sense (and is the actual continuation).

**Weakening neuron 31.9634.** See section I.3.

**Conditional weakening neuron 29.10900.** Gate and linear input weight vectors act as two independent ways of checking that *these* is not present in the residual stream (i.e., a case of double checking). At the same time, they check for predictions like *today, nowadays*. When such predictions are present, the neuron promotes *these*. This is a plausible choice in these cases because of the expression *these days*. An example is *social media tools change and come and go at the drop of a hat*. (This sentence talks about a characteristic of current times, so *these days* would indeed be a plausible continuation.)

**Proportional change neuron 30.10972** predicts the token *when* if activated negatively. This happens if *when* is absent from the residual stream (gate condition) and is proportional to the presence of time-related tokens ( $-w_{\text{in}}$ ). An example for a large negative activation is *puts you on multiple webpages*

<sup>15</sup>The notation is “layer.neuron”, with zero-based indexing.

1188 *at*.<sup>16</sup> Conversely, if *when* is absent, and time-related tokens are absent too, the neuron activates  
1189 positively and suppresses *when* further.

1190 **Orthogonal output neuron 29.4180** predicts *there* (positive activation) if the residual stream contains  
1191 a component that we interpret as “complement of place expected” (e.g., *here*, *therein*). Both  $w_{\text{gate}}$  and  
1192  $w_{\text{in}}$  check for (different aspects of) this component being present, another case of double checking.  
1193 The largest positive activation is on *here or*.  
1194

1195 Overall, these neurons all promote a specific set of tokens (we chose them that way), but under very  
1196 different circumstances. The (conditional) strengthening neurons are the most straightforward to  
1197 interpret, because their input and output clearly correspond to the same concept. In contrast, weaken-  
1198 ing neurons inherently involve (an apparent) conflict between the intermediate model prediction and  
1199 what the neuron promotes.  
1200

## 1201 I.5 RESULTS AND ANALYSIS FOR PROTOTYPICAL EXAMPLES

1202 See table 4.

### 1204 **Strengthening neuron 25.9997.**

1205 **Conditional strengthening neuron 5.10602** activates on those tokens that often start negated auxil-  
1206 iary verbs: *don*, *aren*, *won*, *didn*. Correspondingly, the top token of  $W_U w_{\text{gate}}$  is *t* (but interestingly not  
1207 an apostrophe). On the other hand,  $w_{\text{in}}$  detects alternative predictions: *ate*, *ating* etc. (as in *donate*)  
1208 lead to a negative activation, and *as* (as in *arenas*) leads to a positive activation. Correspondingly  
1209 the strongest positive activations are on the *aren* of *arenas* (but the strongest negative activations are  
1210 not always in a *donate* context, perhaps because both *don't* and *donate* can appear in the same slots).  
1211 These alternative predictions are then strengthened by  $w_{\text{out}}$ .  
1212

### 1213 **Weakening neuron 31.7117.**

### 1214 **Conditional weakening neurons 23.6543.**

### 1215 **Proportional change neuron 25.7415.**

## 1217 I.6 MORE CASE STUDIES

1218 These are various neurons that popped out to us as possibly interesting, for not very systematic  
1219 reasons, for example because they strongly activated on a specific named entity. All of them are in  
1220 OLMo-7B. We present them by RW class. For most of these case studies we did only a quick and  
1221 dirty weight-based analysis. In some cases we also tried  $W_E$  (input embeddings) instead of  $W_U$   
1222 (unembeddings) for the logit-lens style analysis.  
1223  
1224

### 1225 I.6.1 CONDITIONAL STRENGTHENING NEURONS

1226 **0.1480:**  $w_{\text{gate}}$ ,  $-w_{\text{in}}$ ,  $-w_{\text{out}}$  all have tokens similar to *box* (when using  $W_E$ ). Activates on *Xbox*.

1227 **4.1940:** *country* appears in  $w_{\text{in}}$  among many other things. When using  $W_E$ , *Philippines* and *Manila*  
1228 appear in  $w_{\text{out}}$ . Activates on *Philippines*.  
1229

1230 **4.3720:** gate seems country/government related. When using  $W_E$ , we find  $w_{\text{out}}$ ,  $w_{\text{gate}}$  contain some  
1231 country names. Activates on *Denmark*.  
1232

1233 **4.4801:** *Muhammad* appears in the gate vector. Activates on *Muhammad*.  
1234

1235 **4.5772:** predicts *ian* as in *Egyptian*. When using  $W_E$ , all three weight vectors contain *Egypt*. Activates  
1236 on *Egypt*.  
1237

1238 <sup>16</sup>The actual sentence ends with *as soon as* and comes from a now-dead webpage. We  
1239 also found one occurrence of *at when* in what seems to be a paraphrase of the same text, on  
1240 <https://www.docdroid.net/RgxdG5s/fantastic-tips-for-bloggers-of-all-amountsoystepdf-pdf>. We suspect that  
1241 both texts are machine-generated paraphrases of an original text containing *at once* (*when* and *as soon as* can be  
synonyms of *once* in other contexts), and that the model has (also) seen a paraphrased version with *at when*. In  
fact many of the largest negative activations are on *at* in contexts calling for *at once*.

- 1242 **4.6517** has a very Ireland (or Celtic nations) related gate vector. The interpretations of the other two  
 1243 weights are less obvious, but *Irish* and *Dublin* appear in  $w_{in}$  among many other things, and *UK* and  
 1244 *London* appear in  $-w_{out}$  (Ireland is emphatically *not* in the UK!) When using  $W_E$ , *Ireland* appears  
 1245 among the top tokens of all three weight vectors. Activates on *Ireland*.
- 1246 **4.6799**: When using  $W_E$ , *Vietnam* is among the tokens corresponding to  $-w_{out}$ . Activates on *Vietnam*
- 1247 **4.7667**: all three weights related to consoles in different ways. Activates on *Xbox*
- 1249 **4.9983**:  $w_{out}$  is related to electronic devices,  $w_{in}$  either electronic devices or sports (surfing may  
 1250 belong to both),  $w_{gate}$  is also mostly related to electronic devices. When using  $W_E$ , we find  $w_{out}$   
 1251 contains *iPhone* as a top token. Activates on *iPhone*.
- 1252 **4.10859**: When using  $W_E$ , we find  $w_{gate}$ ,  $w_{out}$  include *Thailand* as a top token,  $w_{out}$  additionally  
 1253 *Buddha*, *Buddhist*. Activates on *Thailand*.
- 1254 **4.10882**: When using  $W_E$ , we find  $-w_{out}$  contains *Italy*,  $-w_{in}$ ,  $w_{gate}$  additionally contain *Rome*.  
 1255 Activates on *Italy*.
- 1256 **4.10995**: *Boston* appears in gate and *Massachusetts* in  $-w_{in}$ . When using  $W_E$ , we find  $-w_{out}$ ,  $w_{gate}$   
 1257 contain *Massachusetts* and *Boston*,  $-w_{in}$  contains *Boston*. Activates on *Massachusetts*.
- 1258 **22.2589**:  $w_{gate}$  and  $-w_{in}$  recognize tokens like *Islam*, *Muhammad* and others related to the Arabo-  
 1259 Islamic world. The same goes for  $-w_{out}$  (as it is similar to  $w_{in}$ ). Activates on *Muhammad*.
- 1260 **24.4880**: For all three weight vectors the first four tokens (but not more) are Philippine-related (even  
 1261 though the gate vector is actually not very similar to the others). The gate vector also reacts to other  
 1262 geographical names, which *may* have in common that they are associated with non-”white” (Black,  
 1263 Asian or Latin) people in the US sense (*Singapore*, *Malaysian*, *Nigerian*, *Seoul*, *Pacific*, *Kerala*,  
 1264 *Bangkok*, but also (*Los Angeles* and *Bronx*). Activates on *Philippines*.
- 1265 **24.6771**:  $w_{gate}$ ,  $-w_{in}$ ,  $-w_{out}$  all correspond to capitalized first names. Activates on *Muhammad*.
- 1266 **25.2723**: Some tokens associated with  $w_{in}$  and  $w_{out}$  are possible completions for *th* (*th-ousand*,  
 1267 *th-ought*, *th-orn*). When using  $W_E$ , in all three weights there are a few *th* tokens, but also with *ph* and  
 1268 similar. Activates on *Thailand*.
- 1269 **25.10496**:  $-w_{in}$ ,  $-w_{out}$  correspond to tokens starting with *v* (upper or lower case, with or without  
 1270 preceding space).  $w_{gate}$  on the other hand seems to react to appropriate endings for tokens starting  
 1271 in *v*: *vol-atility*, *v-antage*, *v-intage*, *vel-ocity*, *V-ancouver*. When using  $W_E$ , we also find all three  
 1272 weight-vectors are very *v-heavy*. Activates on *Vietnam*.
- 1273 **I.6.2 WEAKENING NEURONS**
- 1274 **30.9996**: Downgrades weird tokens if present / promotes frequent English stopwords if absent. Also  
 1275 an attention deactivation neuron for 15 heads in layer 31.
- 1276 **I.6.3 PROPORTIONAL CHANGE NEURONS**
- 1277 **25.7032**: Some tokens associated with  $w_{gate}$  and  $w_{out}$  are possible completions for *x* or *ex* (*X-avier*,  
 1278 *x-yz*, *ex-cel*, *ex-ercise*). When using  $W_E$ , both *x* and *box* (with variants) appear in all three weight  
 1279 vectors. Activates on *Xbox*.
- 1280 **25.8607**: All three vectors correspond to tokens related to cities. Moreover,  $-w_{out}$  seems to corre-  
 1281 spond to non-city places, such as national governments or villages.  $w_{in}$  is actually not that similar  
 1282 to  $w_{gate}$ ,  $w_{out}$  (in terms of cosine similarities), but all three correspond to city-related tokens. When  
 1283 using  $W_E$ , in all three weights there are a few city-related tokens. Activates on *Paris*. We may think  
 1284 of the two input directions as two largely independent ways of checking that “it’s about a city” (this  
 1285 is a recurring phenomenon that we describe in section H). When the gate activates but the linear input  
 1286 does not confirm it’s about a city, the output promotes closely related but non-city interpretations (for  
 1287 example *Paris* actually refers to the French government in some contexts).
- 1288 **29.8118**: Partition neuron, highest variance of all proportional change neurons. Also an attention  
 1289 deactivation neuron for 4 heads (0,2,11,15) in layer 30.

- 1296 **31.5490:** Activates on *Muhammad*.  $w_{\text{gate}}$  reacts to various Asian names and Asian-sounding subwords,  
 1297  $w_{\text{in}}$  to surnames as opposed to other English words starting with space and uppercase letter.  $w_{\text{out}}$   
 1298 corresponds to more Asian stuff (mostly subwords) as opposed to English surnames.
- 1299 **31.6275:** Mostly promotes two-letter tokens (no preceding space, typically uppercase).  $-w_{\text{in}}$  typically  
 1300 lowercase single letters.  $-w_{\text{gate}}$  mostly lowercase two-letter tokens. "If no lowercase two-letter  
 1301 tokens, promote uppercase two-letter tokens proportionally to absence of lowercase single letters" ?  
 1302
- 1303 **31.8342:** This is an *-ot-* neuron:  $w_{\text{gate}}$  and  $w_{\text{out}}$  correspond to *-o(t)-* suffixes,  $-w_{\text{in}}$  to various *-ot-* stuff.  
 1304 Judging by the weight similarities, we expect that  $w_{\text{out}}$  is typically activated negatively: downgrade  
 1305 *-o(t)-* suffixes if present in the residual stream. Activates on *Egypt*.
- 1306
- 1307 I.6.4 ORTHOGONAL OUTPUT NEURONS
- 1308 **0.1758:** When using  $W_E$ , all three weight vectors' top tokens are famous web sites, including  
 1309 *YouTube*. Activates on *YouTube*.
- 1310 **0.3338:** When using  $W_E$ , we find especially  $w_{\text{gate}}$  and  $-w_{\text{in}}$ , but also  $-w_{\text{out}}$  are similar to smartphone-  
 1311 related tokens. Activates on *iPhone*.
- 1312 **0.3872:** When using  $W_E$ , we find especially  $w_{\text{gate}}$ , but also  $-w_{\text{in}}$  and  $-w_{\text{out}}$  correspond to city  
 1313 names. Activates on *Paris*.
- 1314 **0.7829:** When using  $W_E$ , we find  $w_{\text{in}}$ ,  $w_{\text{out}}$  and to a lesser extent  $w_{\text{gate}}$  correspond in large part to  
 1315 software names. Activates on *iTunes*.
- 1316 **0.7966:** When using  $W_E$ , the weight vectors mostly correspond to tokens starting with *th*. Activates  
 1317 on *Thor*.
- 1318 **29.2568:**  $w_{\text{out}}$  Asian (Thai?) sounding syllables vs. (Asian) geographic names in English and other  
 1319 stuff;  $w_{\text{in}}$  reacts to Thailand and Asian (geography) stuff as opposed to (mostly) US stuff;  $w_{\text{gate}}$  pretty  
 1320 much the same. Activates on *Thailand*.
- 1321 **29.3327:**  $w_{\text{gate}}$  mostly reacts to city names (*Paris* being the most important one),  $-w_{\text{in}}$  countries  
 1322 and cities, especially in continental Europe (*France* and *Paris* on top) as opposed to stuff related to  
 1323 the former British Empire. Relevant is  $-w_{\text{out}}$  which corresponds to pieces of geographical names  
 1324 and especially rivers in France (*Se-ine*, *Rh-one* / *Rh-ine*, *Mar-ne*, *Mos-elle*... *Norm-andie*, *Nancy*,  
 1325 *commun*...).  $w_{\text{gate}}$  and  $-w_{\text{in}}$  also react to *river(s)*. Activates on *Paris*.
- 1326 **29.4101:**  $w_{\text{gate}}$  and  $w_{\text{in}}$  react to *YouTube* (top token!),  $w_{\text{out}}$  downgrades it (almost bottom token) and  
 1327 promotes *subscrib\**, *views*, *channels* etc. Activates on *YouTube*.
- 1328 **29.6417:** Downgrades *recording* and similar.  $w_{\text{gate}}$  and  $w_{\text{in}}$  are also similar and involve *iTunes*.  
 1329 Activates on *iTunes*.
- 1330 **29.9734:**  $w_{\text{gate}}$  reacts to the East in a broad sense as opposed to the West (*Iran*, *Kaz-akhstan*, *Kash-mir*,  
 1331 *Ukraine*...),  $w_{\text{in}}$  mostly to male first names without preceding space.  $w_{\text{out}}$  seems to produce word  
 1332 pieces that could begin a foreign name. Activates on *Muhammad*.
- 1333 **30.2667:**  $w_{\text{gate}}$  reacts to suffixes (for adjectives derived from place names) like *en*, *ian*, *ians*, basically  
 1334 the same for  $w_{\text{in}}$  and  $w_{\text{out}}$ . Activates on *Muhammad*.
- 1335 **30.3143:**  $w_{\text{gate}}$  reacts to words related to entities that are authoritative for various reasons (*officials*,  
 1336 *authorities*, *according*, *researchers*, *spokesman*, *investigators*...).  $-w_{\text{in}}$  reacts to uncertainty (*report-*  
 1337 *edly*, *according*... *allegedly*... *accused*).  $-w_{\text{out}}$  is again *police*, *authorities*, *officials*, *court* but with  
 1338 no preceding space. Activates on *Philippines*. What authorities and uncertainty have to do with the  
 1339 Philippines is unclear.
- 1340 **30.3883:**  $w_{\text{gate}}$  and  $-w_{\text{in}}$  react to *Virginia* and *Afghanistan*, among others (in the case of  $w_{\text{gate}}$ : as  
 1341 opposed to other geographical names with no preceding space associated with the South and the sea);  
 1342  $-w_{\text{out}}$  is activated and promotes all variants of *af* (and *ghan*) but downgrades *Virginia* etc. Activates  
 1343 on *Afghanistan*.
- 1344 **30.4577:** Seems to be related to rugby:  $w_{\text{gate}}$  and slightly less obviously  $w_{\text{in}}$  react to rugby-related  
 1345 tokens (*midfielder*, *quarterback*...);  $w_{\text{out}}$  promotes different tokens that upon reflection could be  
 1346 related to rugby as well. Activates on *Ireland*.

1350 **30.5372:** Promotes *natural* and related, downgrades *inst* tokens.  $w_{in}$  reacts to *wildlife* etc. as opposed  
 1351 to *institute* etc,  $w_{gate}$  reacts to *institute* as opposed to *natural*. Activates on *Massachusetts* (in which  
 1352 situation it promotes *Institute*, which makes sense because of MIT).

1353 **30.8535:**  $-w_{out}$  is *one* in all variants,  $w_{gate}$  too,  $w_{in}$  splits *one*, *ones* and the equivalent Chinese  
 1354 characters, on the positive side, from *One*, *I*, *ONE* on the negative side (and many other things on  
 1355 both sides). Activates on *Xbox*. Presumably this happens because *One* is a possible prediction (*Xbox*  
 1356 *One*), and presumably the output reinforces that.

1357 **31.2135:** orthogonal output, on the conditional strengthening side (weak conditional strengthening,  
 1358 one of the neurons on the vertical axis).  $w_{gate}$  reacts to single letters or symbols as opposed to some  
 1359 English content words without preceding space;  $w_{in}$  and  $w_{out}$  mostly Chinese or Japanese characters  
 1360 as opposed to some Latin diacritics and other weird stuff. Language choice? "If it's not English and  
 1361 single letters are floating around, make sure to choose the right language / character set."

1362 **31.10424:**  $w_{gate}$ ,  $-w_{in}$ ,  $w_{out}$  correspond to *score* in the top tokens, which is downgraded if present.  
 1363 Activates on *Paris*. No idea what's happening here.

## 1364 J MORE PLOTS

1365  
 1366 These final figures show additional results:

- 1367
- 1368
- 1369
- 1370 • Table 5 and figures 7 to 9 show more results on activation frequencies: by discrete classes,  
 1371 on all layers, and plotted against  $\cos(w_{gate}, w_{in})$  or  $\cos(w_{gate}, w_{out})$ .
- 1372 • Figures 14 to 16 show that ablating 243 neurons from other classes than weakening does not  
 1373 have any effect on attribute rate.
- 1374 • Figures 17 to 24 show the effect of various ablations on entropy, loss, rank and scale.
- 1375 • From figure 40, we show our analyses of RW functionalities by layer (section 5) for all the  
 1376 models we investigated.
- 1377

1378 Regarding the last point, we note a few additional patterns that appear only in some of these models:

- 1379
- 1380 • In Yi and the OLMo models, the prevalence of conditional strengthening neurons starts  
 1381 even earlier, at the very first layer. A particularly interesting example is Yi: In layer 0 an  
 1382 enormous 68% of all neurons are conditional strengthening, then almost none, then there  
 1383 is a second wave around layers 11-17 (out of 32) which have around 25% of conditional  
 1384 strengthening neurons each.
- 1385 • In some models, especially the OLMo ones, there is a non-negligible number of conditional  
 1386 weakening neurons. They tend to appear in middle-to-late layers, shortly after the conditional  
 1387 strengthening wave. The clearest example is OLMo-1B, with a peak of 1418 conditional  
 1388 weakening neurons out of 8192 (17%) in layer 9 out of 16.

1389 The following patterns could be random, but still show that the model has *not* learned something:

- 1390
- 1391 • For almost all neurons the cosine similarities are still clearly below 1 (the dots do not fill out  
 1392 the edges in figure 2). This echoes and extends Gurnee et al. (2024)'s findings that in GPT2  
 1393 the IO cosine similarity is approximately bounded by  $\pm 0.8$ . In other words, we almost  
 1394 never get the *prototypical* cases of conditional strengthening / weakening etc., as defined in  
 1395 section 4. This might be an effect of randomness (strong cosine similarities are less likely),  
 1396 but could also suggest that even input manipulator neurons add some novel information to  
 1397 the residual stream.
- 1398 • We also observe that for the vast majority of neurons,  $\cos(w_{gate}, w_{in}) \approx 0$ : This can be seen  
 1399 in the boxplots in the appendix, as well as the purple color in figure 2. Thus most neurons  
 1400 operate on two input directions in the residual stream (not a single one), resulting in higher  
 1401 expressivity and more complex semantics. If not random, this could be related to double  
 1402 checking; see section H.
- 1403

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457

Table 5: Activation frequencies by RW class. The second column represents random neurons taken from the same layers.

	true	baseline
strengthening	0.265	0.480
atypical strengthening	0.163	0.276
conditional strengthening	0.132	0.370
atypical conditional strengthening	0.100	0.219
proportional change	0.368	0.296
atypical proportional change	0.344	0.338
orthogonal output	0.373	0.236
weakening	0.691	0.384
atypical weakening	0.727	0.353
conditional weakening	0.708	0.301
atypical conditional weakening	0.665	0.469

1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

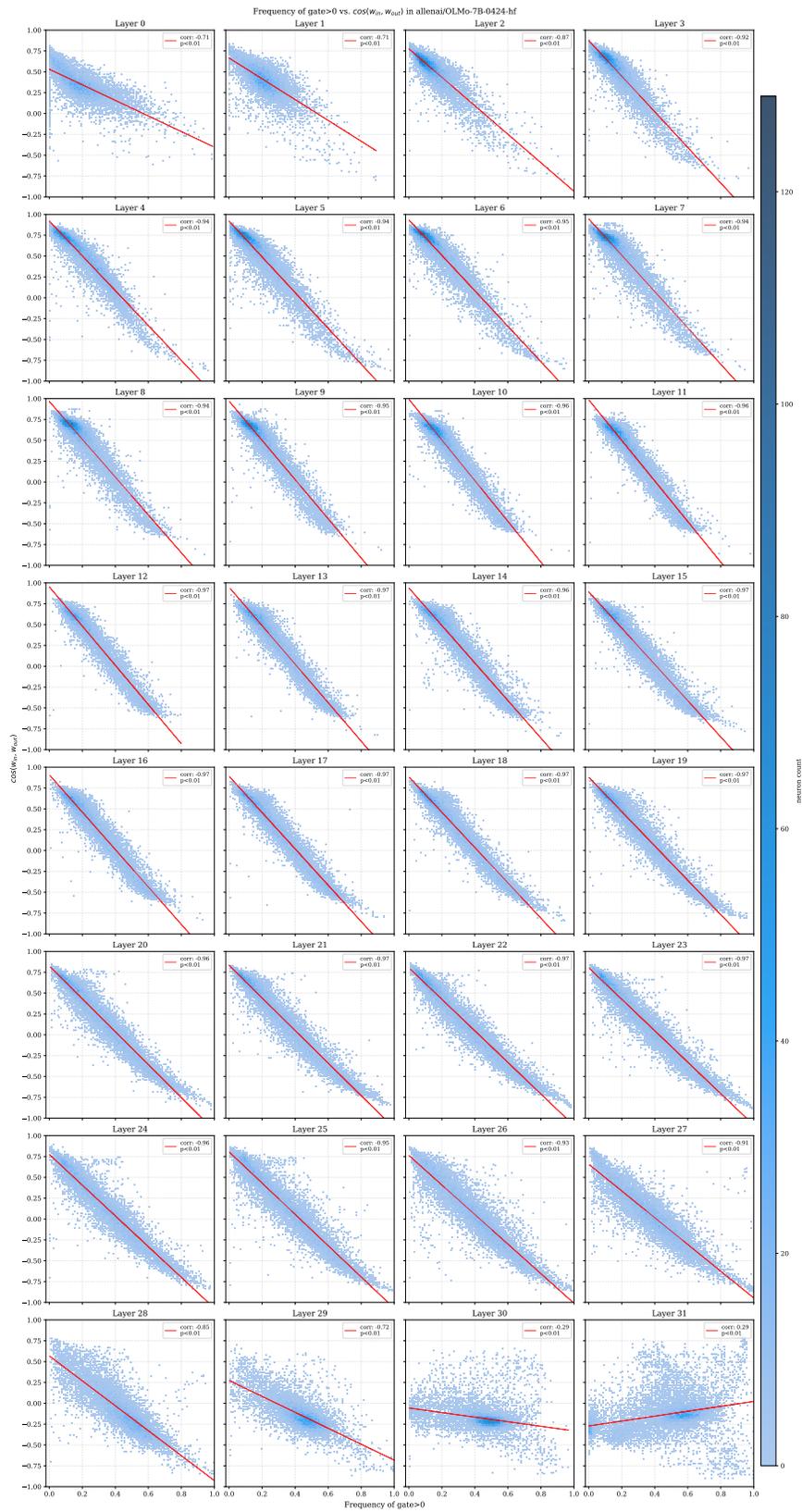


Figure 7: Like figure 4 but for all layers.

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

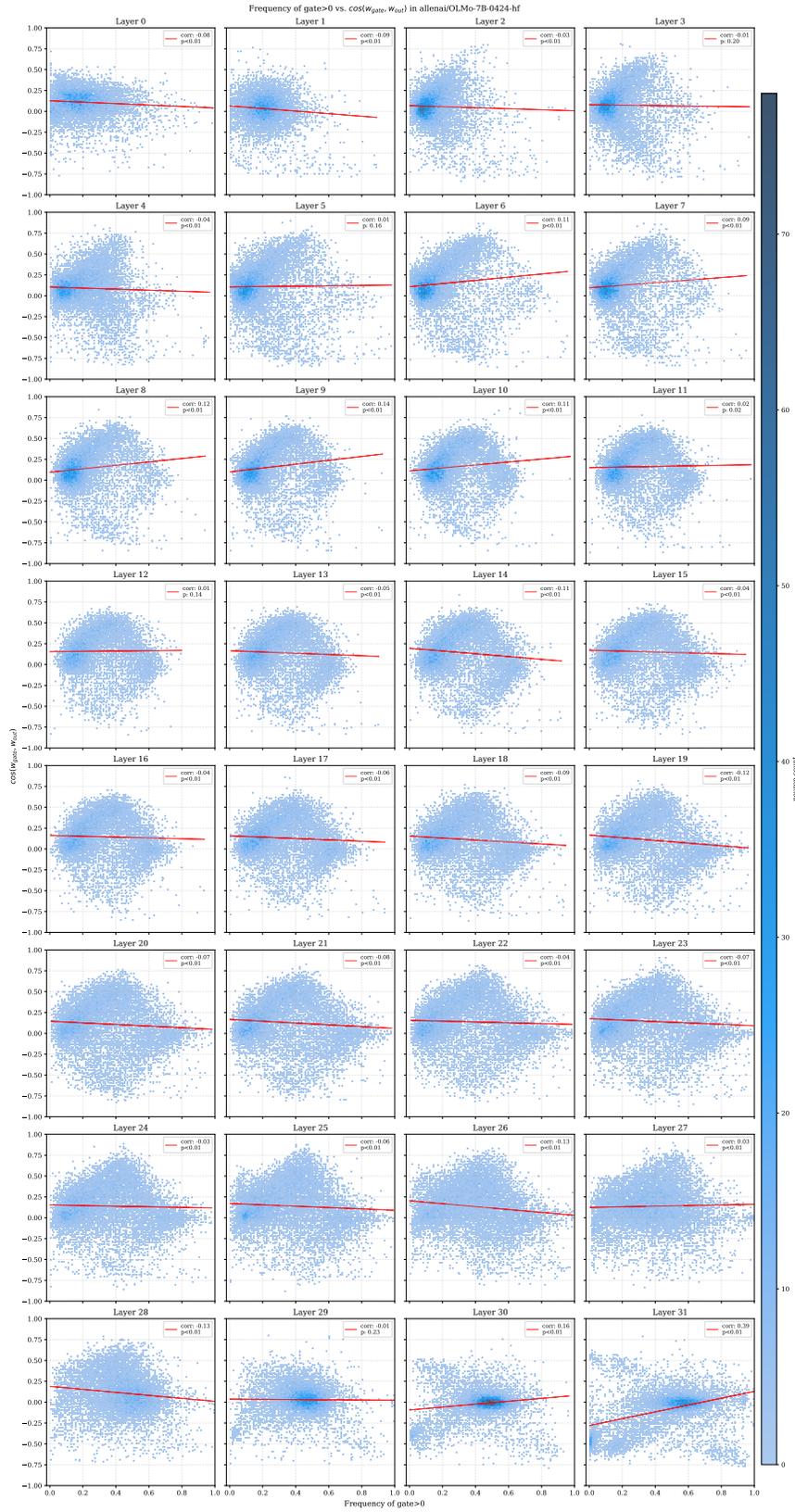


Figure 8: Frequency of gate>0 vs.  $|\cos(w_{\text{gate}}, w_{\text{out}})|$  in OLMo-7B.

1566  
 1567  
 1568  
 1569  
 1570  
 1571  
 1572  
 1573  
 1574  
 1575  
 1576  
 1577  
 1578  
 1579  
 1580  
 1581  
 1582  
 1583  
 1584  
 1585  
 1586  
 1587  
 1588  
 1589  
 1590  
 1591  
 1592  
 1593  
 1594  
 1595  
 1596  
 1597  
 1598  
 1599  
 1600  
 1601  
 1602  
 1603  
 1604  
 1605  
 1606  
 1607  
 1608  
 1609  
 1610  
 1611  
 1612  
 1613  
 1614  
 1615  
 1616  
 1617  
 1618  
 1619

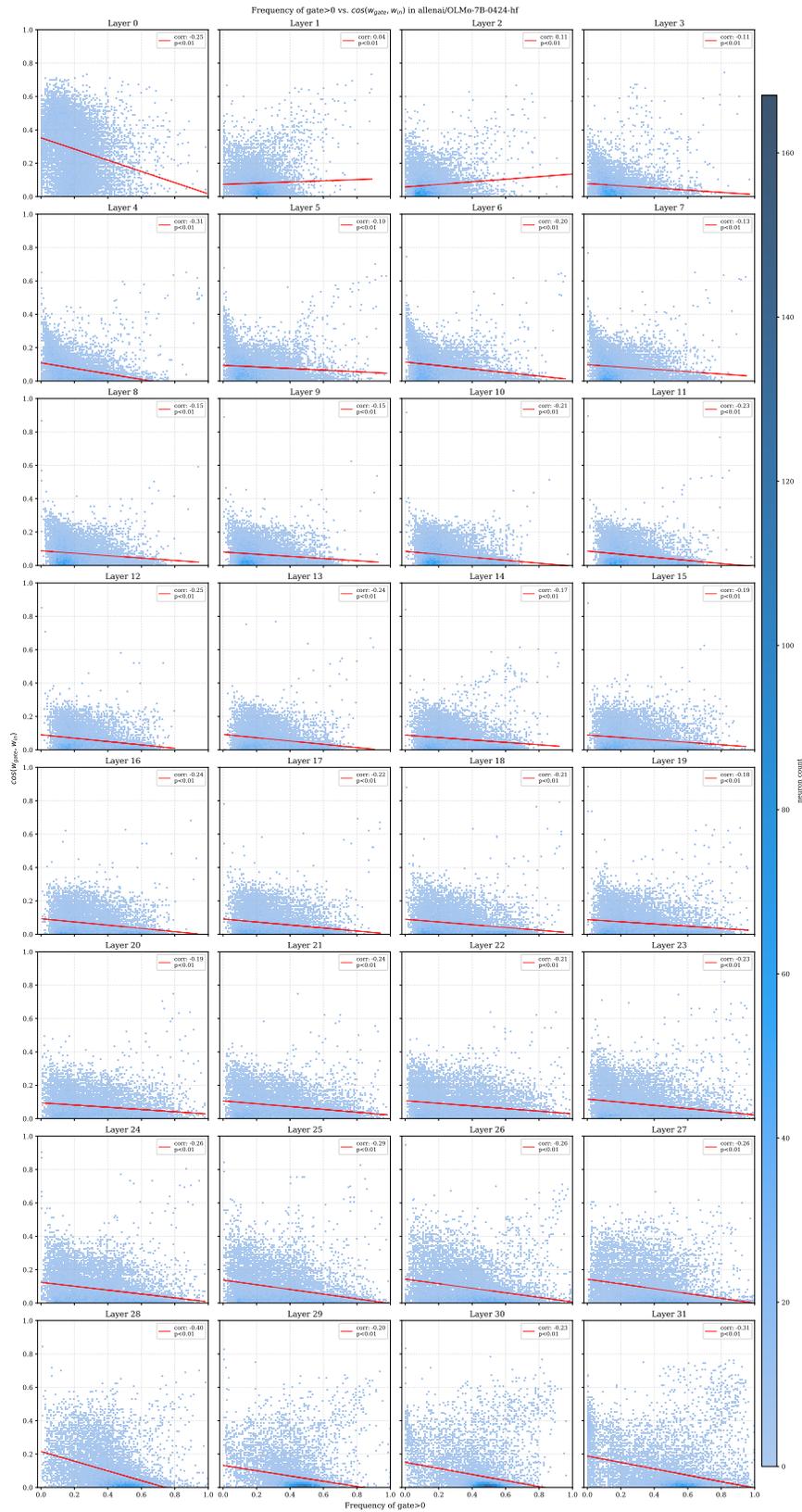


Figure 9: Frequency of gate>0 vs.  $\cos(w_{\text{gate}}, w_{\text{in}})$  in OLMo-7B.

1620  
 1621  
 1622  
 1623  
 1624  
 1625  
 1626  
 1627  
 1628  
 1629  
 1630  
 1631  
 1632  
 1633  
 1634  
 1635  
 1636  
 1637  
 1638  
 1639  
 1640  
 1641  
 1642  
 1643  
 1644  
 1645  
 1646  
 1647  
 1648  
 1649  
 1650  
 1651  
 1652  
 1653  
 1654  
 1655  
 1656  
 1657  
 1658  
 1659  
 1660  
 1661  
 1662  
 1663  
 1664  
 1665  
 1666  
 1667  
 1668  
 1669  
 1670  
 1671  
 1672  
 1673

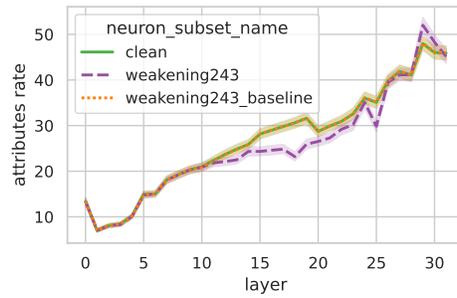


Figure 10: Attribute rate when mean-ablating all 243 weakening neurons.

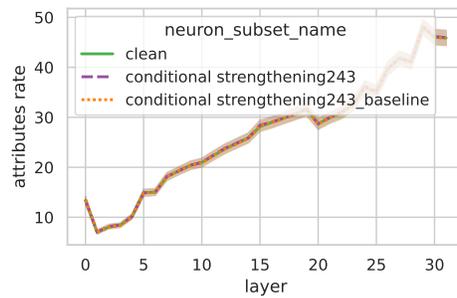


Figure 11: Attribute rate when mean-ablating 243 conditional strengthening neurons.

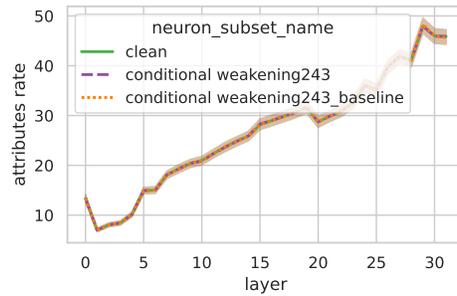


Figure 12: Attribute rate when mean-ablating 243 conditional weakening neurons.

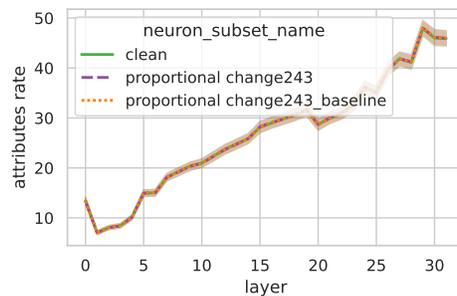


Figure 13: Attribute rate when mean-ablating 243 proportional change neurons.

1674  
 1675  
 1676  
 1677  
 1678  
 1679  
 1680  
 1681  
 1682  
 1683  
 1684  
 1685  
 1686  
 1687  
 1688  
 1689  
 1690  
 1691  
 1692  
 1693  
 1694  
 1695  
 1696  
 1697  
 1698  
 1699  
 1700  
 1701  
 1702  
 1703  
 1704  
 1705  
 1706  
 1707  
 1708  
 1709  
 1710  
 1711  
 1712  
 1713  
 1714  
 1715  
 1716  
 1717  
 1718  
 1719  
 1720  
 1721  
 1722  
 1723  
 1724  
 1725  
 1726  
 1727

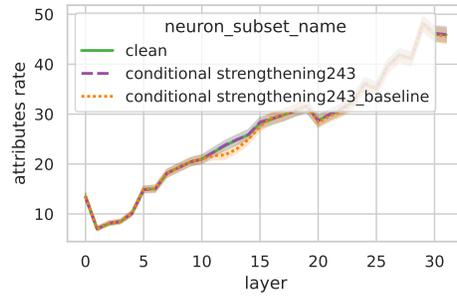


Figure 14: Attribute rate when zero-ablating 243 conditional strengthening neurons.

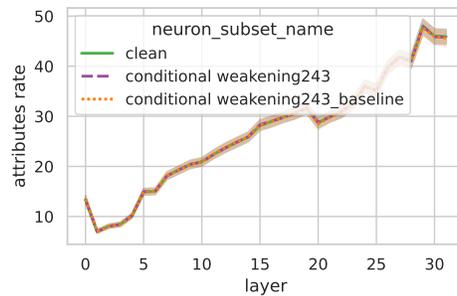


Figure 15: Attribute rate when zero-ablating 243 conditional weakening neurons.

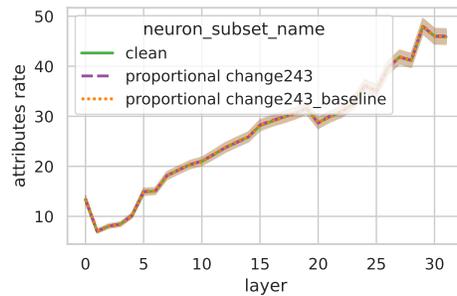


Figure 16: Attribute rate when zero-ablating 243 proportional change neurons.

1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781

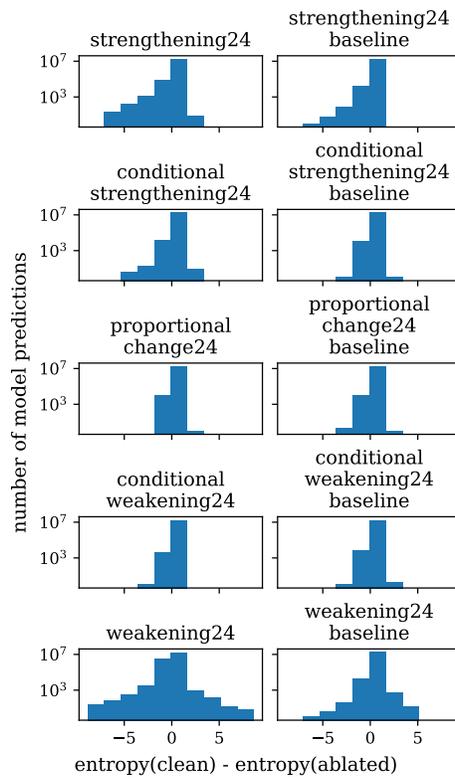


Figure 17: Effect on entropy when zero-ablating 24 neurons from various RW classes.

1782  
 1783  
 1784  
 1785  
 1786  
 1787  
 1788  
 1789  
 1790  
 1791  
 1792  
 1793  
 1794  
 1795  
 1796  
 1797  
 1798  
 1799  
 1800  
 1801  
 1802  
 1803  
 1804  
 1805  
 1806  
 1807  
 1808  
 1809  
 1810  
 1811  
 1812  
 1813  
 1814  
 1815  
 1816  
 1817  
 1818  
 1819  
 1820  
 1821  
 1822  
 1823  
 1824  
 1825  
 1826  
 1827  
 1828  
 1829  
 1830  
 1831  
 1832  
 1833  
 1834  
 1835

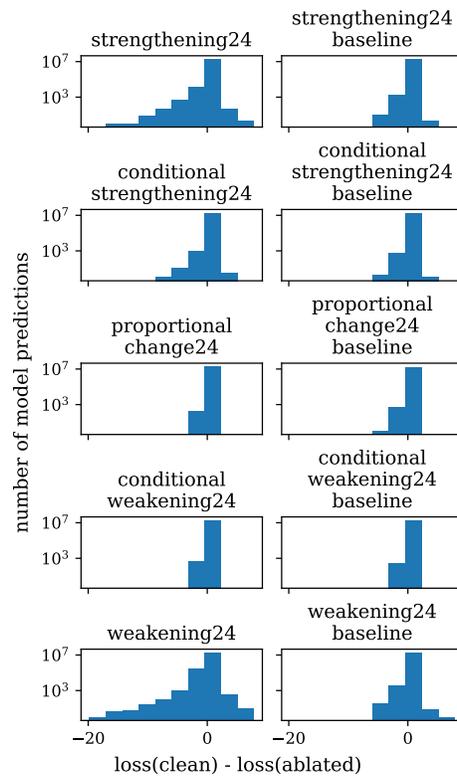


Figure 18: Effect on loss when zero-ablating 24 neurons from various RW classes.

1836  
 1837  
 1838  
 1839  
 1840  
 1841  
 1842  
 1843  
 1844  
 1845  
 1846  
 1847  
 1848  
 1849  
 1850  
 1851  
 1852  
 1853  
 1854  
 1855  
 1856  
 1857  
 1858  
 1859  
 1860  
 1861  
 1862  
 1863  
 1864  
 1865  
 1866  
 1867  
 1868  
 1869  
 1870  
 1871  
 1872  
 1873  
 1874  
 1875  
 1876  
 1877  
 1878  
 1879  
 1880  
 1881  
 1882  
 1883  
 1884  
 1885  
 1886  
 1887  
 1888  
 1889

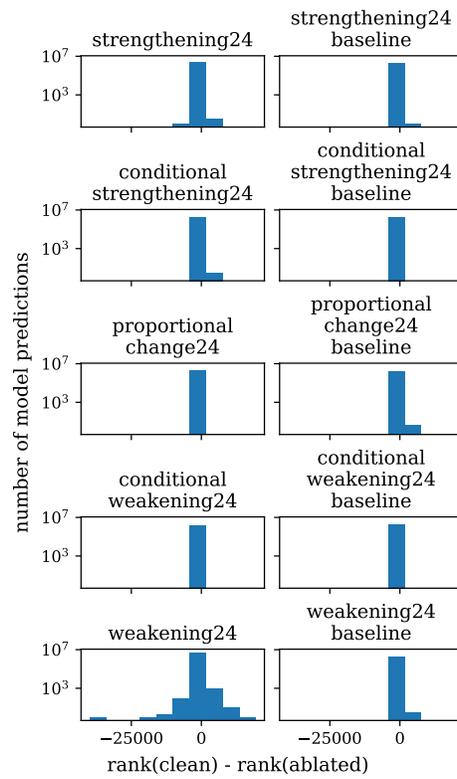


Figure 19: Effect on correct token rank when zero-ablating 24 neurons from various RW classes.

1890  
 1891  
 1892  
 1893  
 1894  
 1895  
 1896  
 1897  
 1898  
 1899  
 1900  
 1901  
 1902  
 1903  
 1904  
 1905  
 1906  
 1907  
 1908  
 1909  
 1910  
 1911  
 1912  
 1913  
 1914  
 1915  
 1916  
 1917  
 1918  
 1919  
 1920  
 1921  
 1922  
 1923  
 1924  
 1925  
 1926  
 1927  
 1928  
 1929  
 1930  
 1931  
 1932  
 1933  
 1934  
 1935  
 1936  
 1937  
 1938  
 1939  
 1940  
 1941  
 1942  
 1943

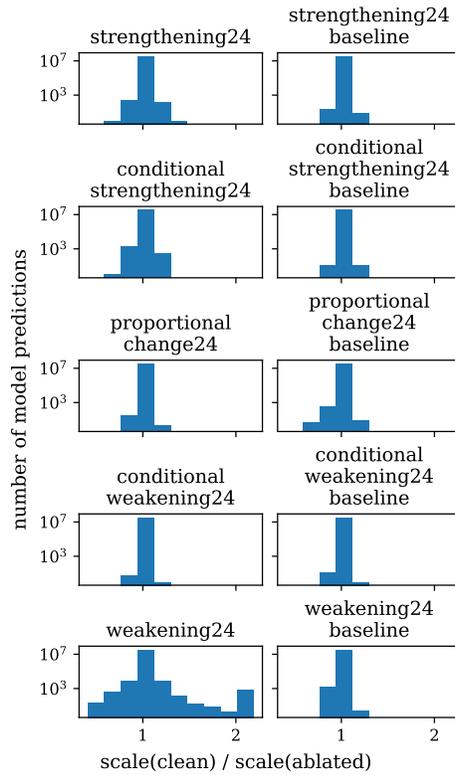


Figure 20: Effect on scale of last hidden state when zero-ablating 24 neurons from various RW classes.

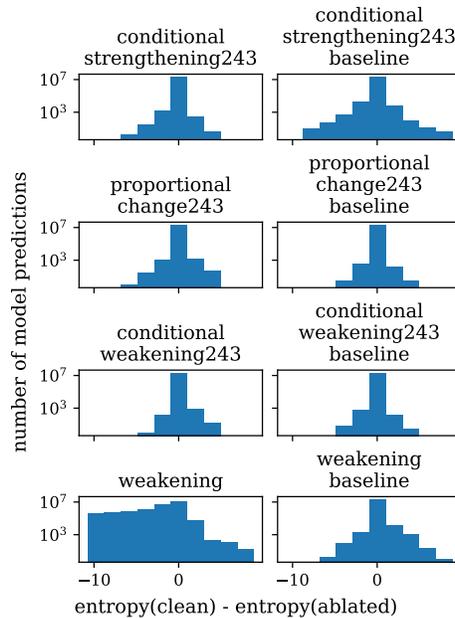


Figure 21: Effect on entropy when zero-ablating 243 neurons from various RW classes.

1944  
 1945  
 1946  
 1947  
 1948  
 1949  
 1950  
 1951  
 1952  
 1953  
 1954  
 1955  
 1956  
 1957  
 1958  
 1959  
 1960  
 1961  
 1962  
 1963  
 1964  
 1965  
 1966  
 1967  
 1968  
 1969  
 1970  
 1971  
 1972  
 1973  
 1974  
 1975  
 1976  
 1977  
 1978  
 1979  
 1980  
 1981  
 1982  
 1983  
 1984  
 1985  
 1986  
 1987  
 1988  
 1989  
 1990  
 1991  
 1992  
 1993  
 1994  
 1995  
 1996  
 1997

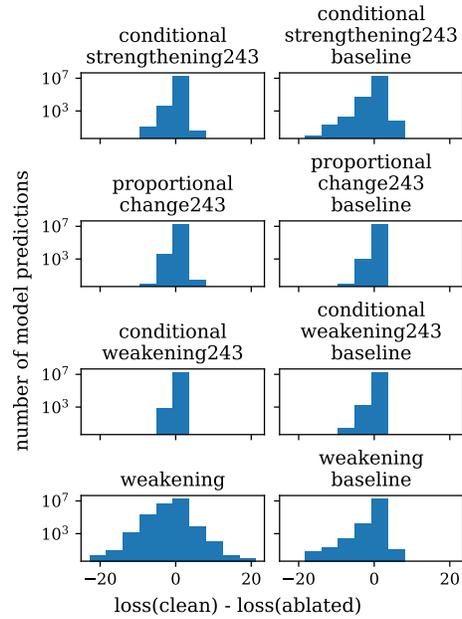


Figure 22: Effect on loss when zero-ablating 243 neurons from various RW classes.

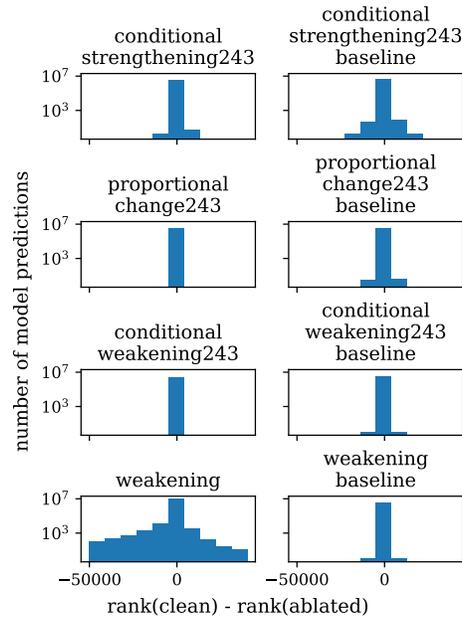


Figure 23: Effect on correct token rank when zero-ablating 243 neurons from various RW classes.

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

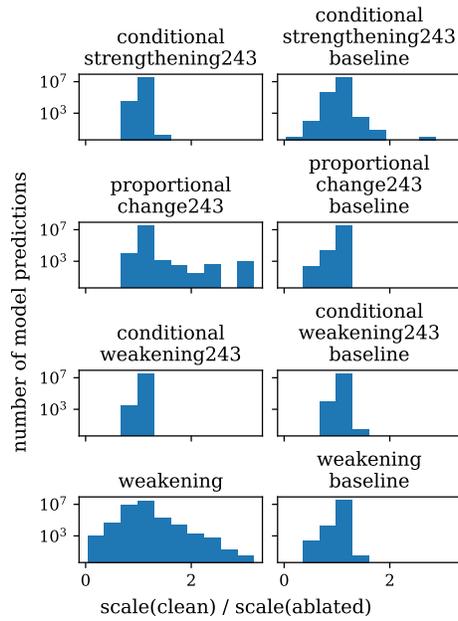


Figure 24: Effect on scale of last hidden state when zero-ablating 243 neurons from various RW classes.

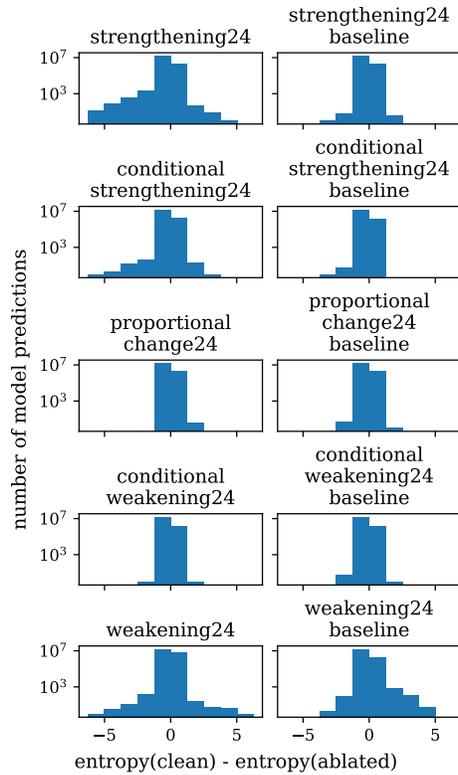


Figure 25: Effect on entropy when mean-ablating 24 neurons from various RW classes.

2052  
 2053  
 2054  
 2055  
 2056  
 2057  
 2058  
 2059  
 2060  
 2061  
 2062  
 2063  
 2064  
 2065  
 2066  
 2067  
 2068  
 2069  
 2070  
 2071  
 2072  
 2073  
 2074  
 2075  
 2076  
 2077  
 2078  
 2079  
 2080  
 2081  
 2082  
 2083  
 2084  
 2085  
 2086  
 2087  
 2088  
 2089  
 2090  
 2091  
 2092  
 2093  
 2094  
 2095  
 2096  
 2097  
 2098  
 2099  
 2100  
 2101  
 2102  
 2103  
 2104  
 2105

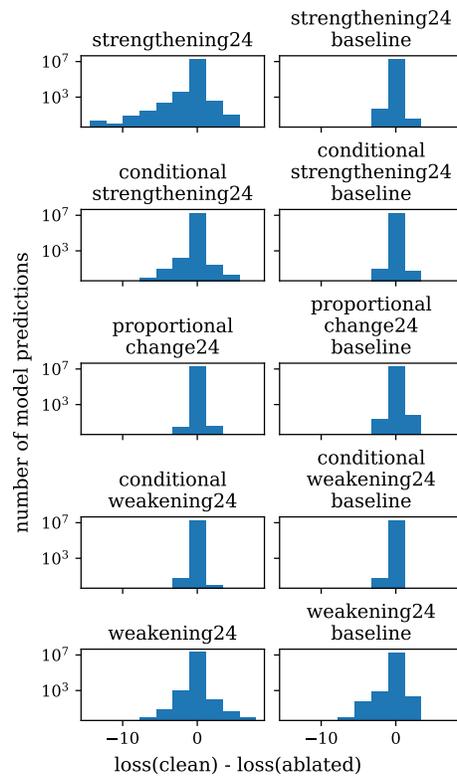


Figure 26: Effect on loss when mean-ablating 24 neurons from various RW classes.

2106  
 2107  
 2108  
 2109  
 2110  
 2111  
 2112  
 2113  
 2114  
 2115  
 2116  
 2117  
 2118  
 2119  
 2120  
 2121  
 2122  
 2123  
 2124  
 2125  
 2126  
 2127  
 2128  
 2129  
 2130  
 2131  
 2132  
 2133  
 2134  
 2135  
 2136  
 2137  
 2138  
 2139  
 2140  
 2141  
 2142  
 2143  
 2144  
 2145  
 2146  
 2147  
 2148  
 2149  
 2150  
 2151  
 2152  
 2153  
 2154  
 2155  
 2156  
 2157  
 2158  
 2159

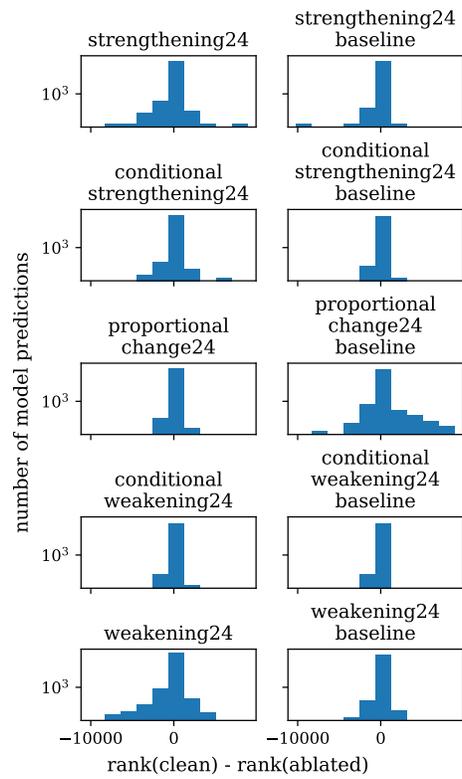
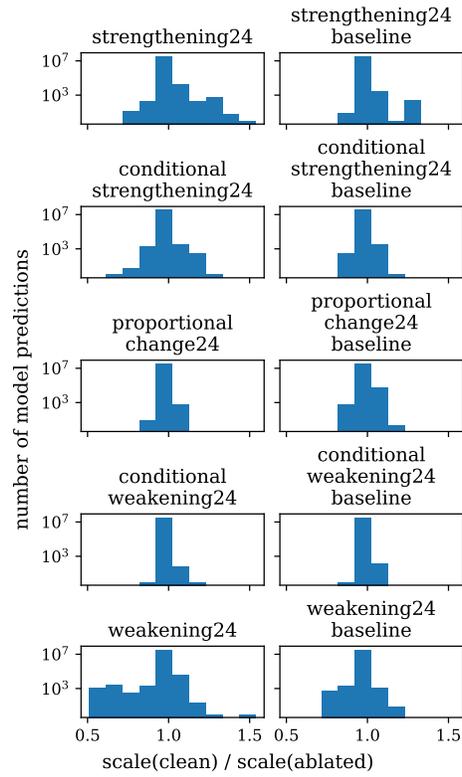


Figure 27: Effect on correct token rank when mean-ablating 24 neurons from various RW classes.

2160  
2161  
2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186



2187 Figure 28: Effect on scale of last hidden state when mean-ablating 24 neurons from various RW  
2188 classes.

2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213

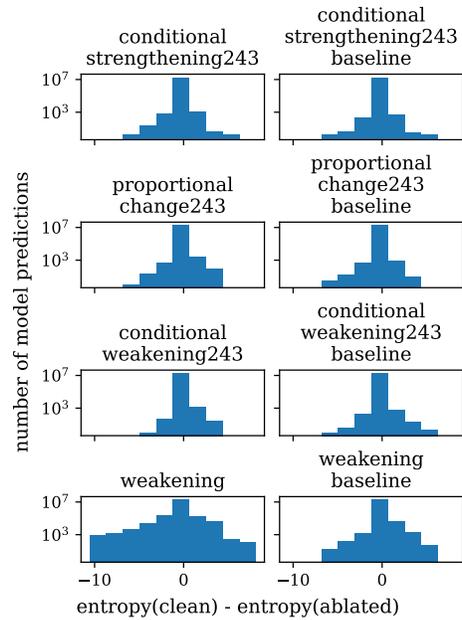


Figure 29: Effect on entropy when mean-ablating 243 neurons from various RW classes.

2214  
 2215  
 2216  
 2217  
 2218  
 2219  
 2220  
 2221  
 2222  
 2223  
 2224  
 2225  
 2226  
 2227  
 2228  
 2229  
 2230  
 2231  
 2232  
 2233  
 2234  
 2235  
 2236  
 2237  
 2238  
 2239  
 2240  
 2241  
 2242  
 2243  
 2244  
 2245  
 2246  
 2247  
 2248  
 2249  
 2250  
 2251  
 2252  
 2253  
 2254  
 2255  
 2256  
 2257  
 2258  
 2259  
 2260  
 2261  
 2262  
 2263  
 2264  
 2265  
 2266  
 2267

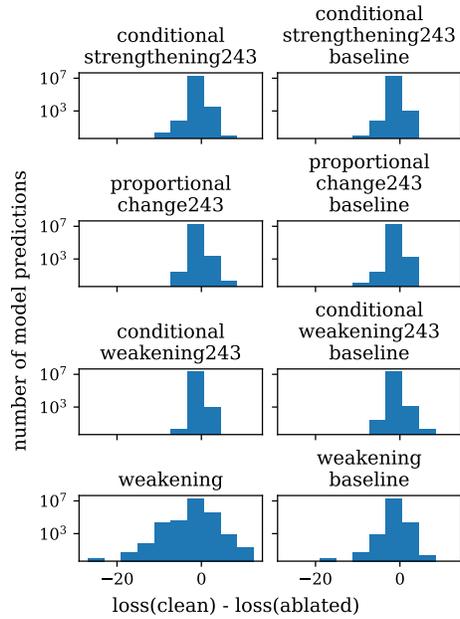


Figure 30: Effect on loss when mean-ablating 243 neurons from various RW classes.

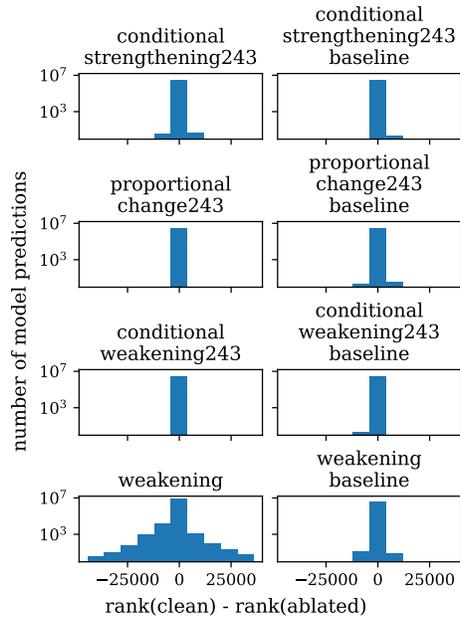


Figure 31: Effect on correct token rank when mean-ablating 243 neurons from various RW classes.

2268  
 2269  
 2270  
 2271  
 2272  
 2273  
 2274  
 2275  
 2276  
 2277  
 2278  
 2279  
 2280  
 2281  
 2282  
 2283  
 2284  
 2285  
 2286  
 2287  
 2288  
 2289  
 2290  
 2291  
 2292  
 2293  
 2294  
 2295  
 2296  
 2297  
 2298  
 2299  
 2300  
 2301  
 2302  
 2303  
 2304  
 2305  
 2306  
 2307  
 2308  
 2309  
 2310  
 2311  
 2312  
 2313  
 2314  
 2315  
 2316  
 2317  
 2318  
 2319  
 2320  
 2321

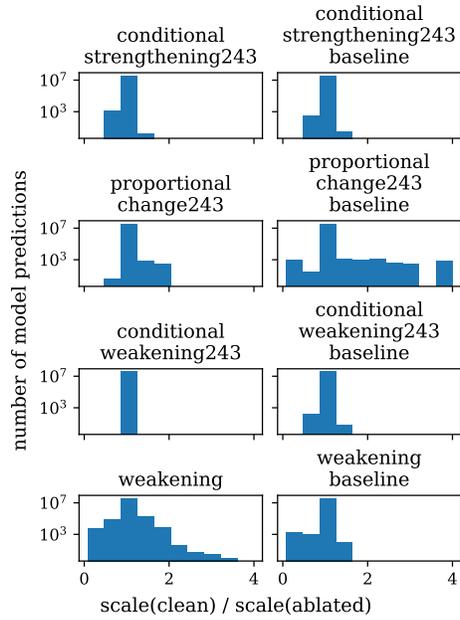


Figure 32: Effect on scale of last hidden state when mean-ablating 243 neurons from various RW classes.

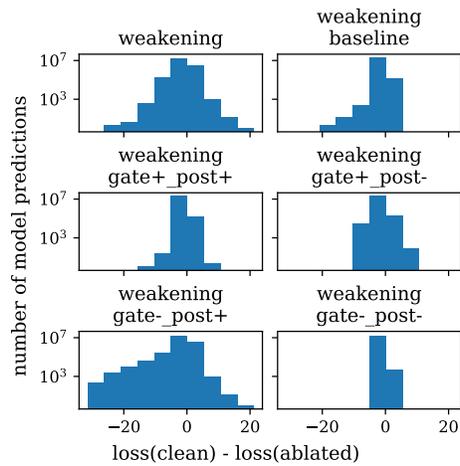


Figure 33: Effect on loss of conditional zero-ablations of weakening neurons.

2322  
2323  
2324  
2325  
2326  
2327  
2328  
2329  
2330  
2331  
2332  
2333  
2334  
2335  
2336  
2337

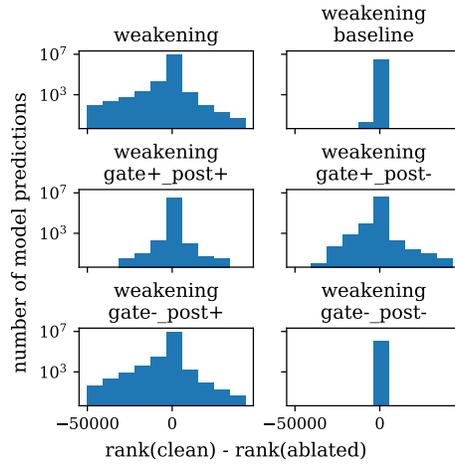


Figure 34: Effect on rank of conditional zero-ablations of weakening neurons.

2338  
2339  
2340  
2341  
2342  
2343  
2344  
2345  
2346  
2347  
2348  
2349  
2350  
2351  
2352  
2353  
2354  
2355  
2356

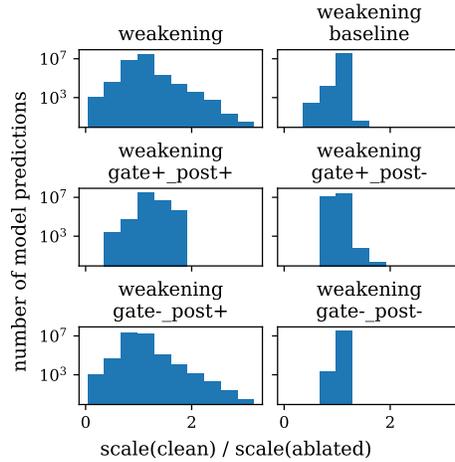


Figure 35: Effect on scale of conditional zero-ablations of weakening neurons.

2357  
2358  
2359  
2360  
2361  
2362  
2363  
2364  
2365  
2366  
2367  
2368  
2369  
2370  
2371  
2372  
2373  
2374  
2375

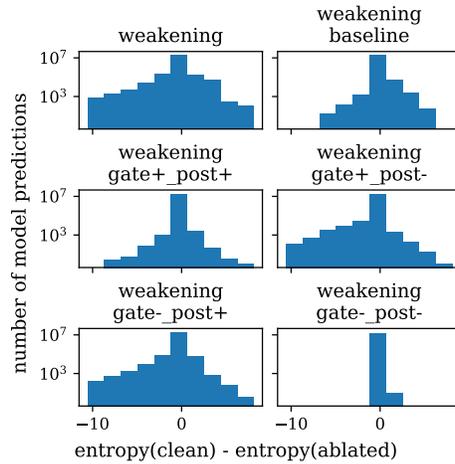


Figure 36: Effect on entropy of conditional mean-ablations of weakening neurons.

2376  
2377  
2378  
2379  
2380  
2381  
2382  
2383  
2384  
2385  
2386  
2387  
2388  
2389  
2390  
2391  
2392  
2393  
2394  
2395  
2396  
2397  
2398  
2399  
2400  
2401  
2402  
2403  
2404  
2405  
2406  
2407  
2408  
2409  
2410  
2411  
2412  
2413  
2414  
2415  
2416  
2417  
2418  
2419  
2420  
2421  
2422  
2423  
2424  
2425  
2426  
2427  
2428  
2429

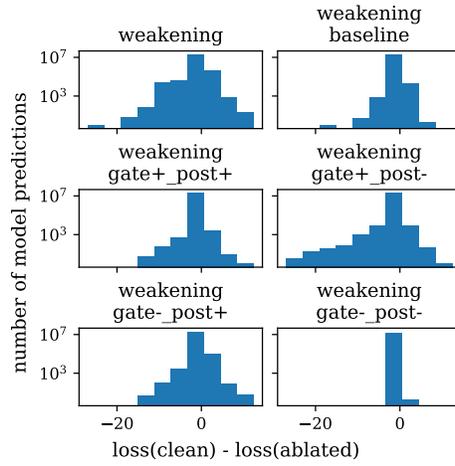


Figure 37: Effect on loss of conditional mean-ablations of weakening neurons.

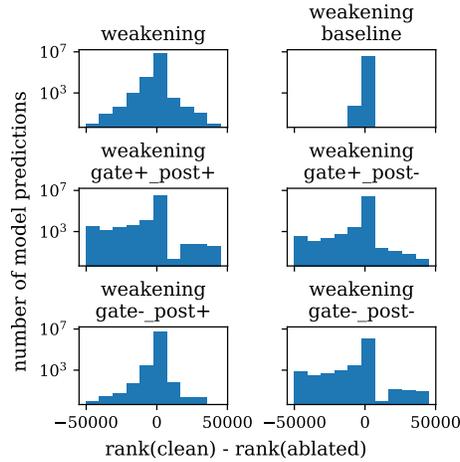


Figure 38: Effect on rank of conditional mean-ablations of weakening neurons.

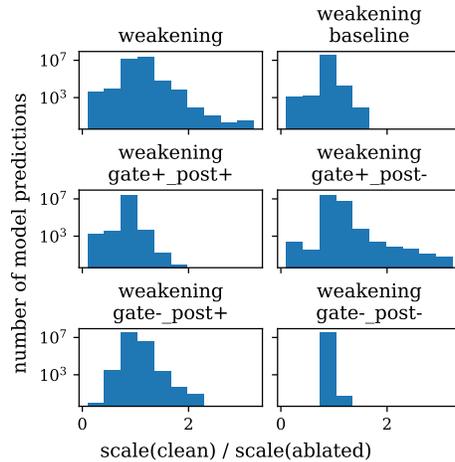


Figure 39: Effect on scale of conditional mean-ablations of weakening neurons.

2430  
 2431  
 2432  
 2433  
 2434  
 2435  
 2436  
 2437  
 2438  
 2439  
 2440  
 2441  
 2442  
 2443  
 2444  
 2445  
 2446  
 2447  
 2448  
 2449  
 2450  
 2451  
 2452  
 2453  
 2454  
 2455  
 2456  
 2457  
 2458  
 2459  
 2460  
 2461  
 2462  
 2463  
 2464  
 2465  
 2466  
 2467  
 2468  
 2469  
 2470  
 2471  
 2472  
 2473  
 2474  
 2475  
 2476  
 2477  
 2478  
 2479  
 2480  
 2481  
 2482  
 2483

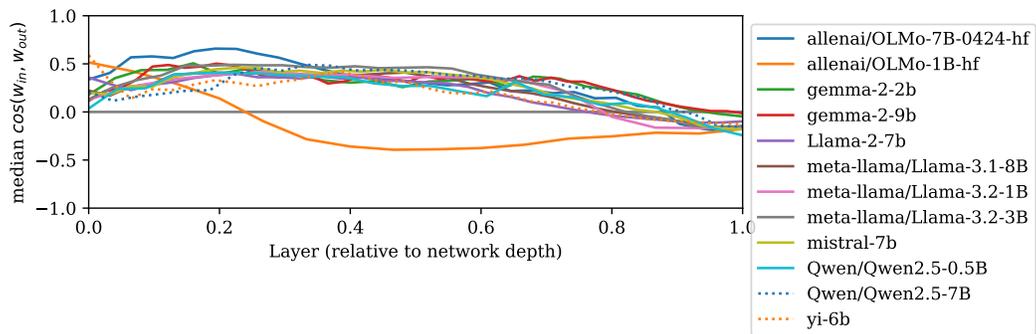


Figure 40: Median of  $\cos(\mathbf{w}_{in}, \mathbf{w}_{out})$  by layer (x-axis) for all 12 models investigated. Unlike figure 1(a) we also include the models of 1B parameters and below. All models follow the same general pattern, but OLMo-1B switches to negative values earlier than the others.

2484  
2485  
2486  
2487  
2488  
2489  
2490  
2491  
2492  
2493  
2494  
2495  
2496  
2497  
2498  
2499  
2500  
2501  
2502  
2503  
2504  
2505  
2506  
2507  
2508  
2509  
2510  
2511  
2512  
2513  
2514  
2515  
2516  
2517  
2518  
2519  
2520  
2521  
2522  
2523  
2524  
2525  
2526  
2527  
2528  
2529  
2530  
2531  
2532  
2533  
2534  
2535  
2536  
2537

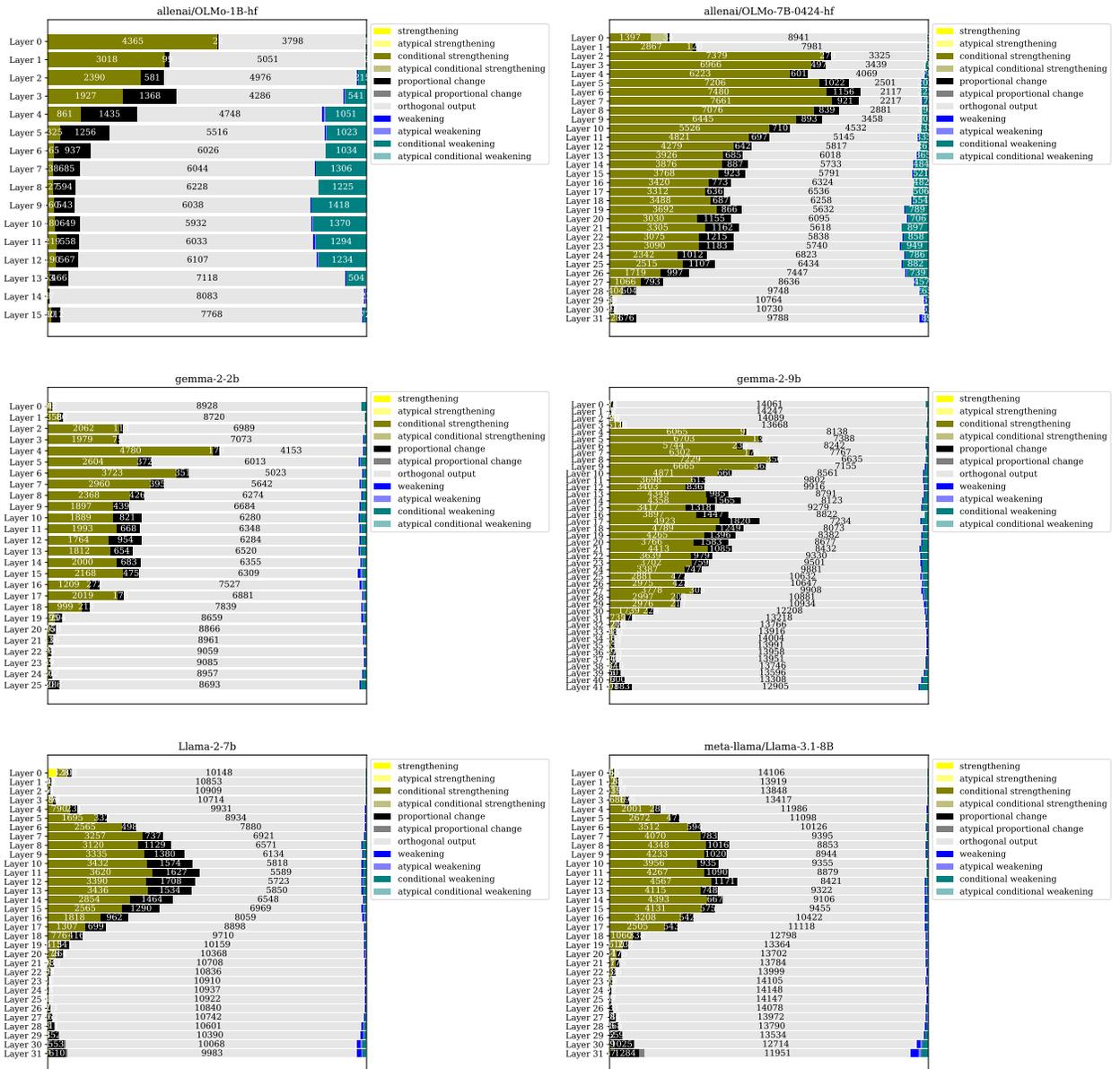
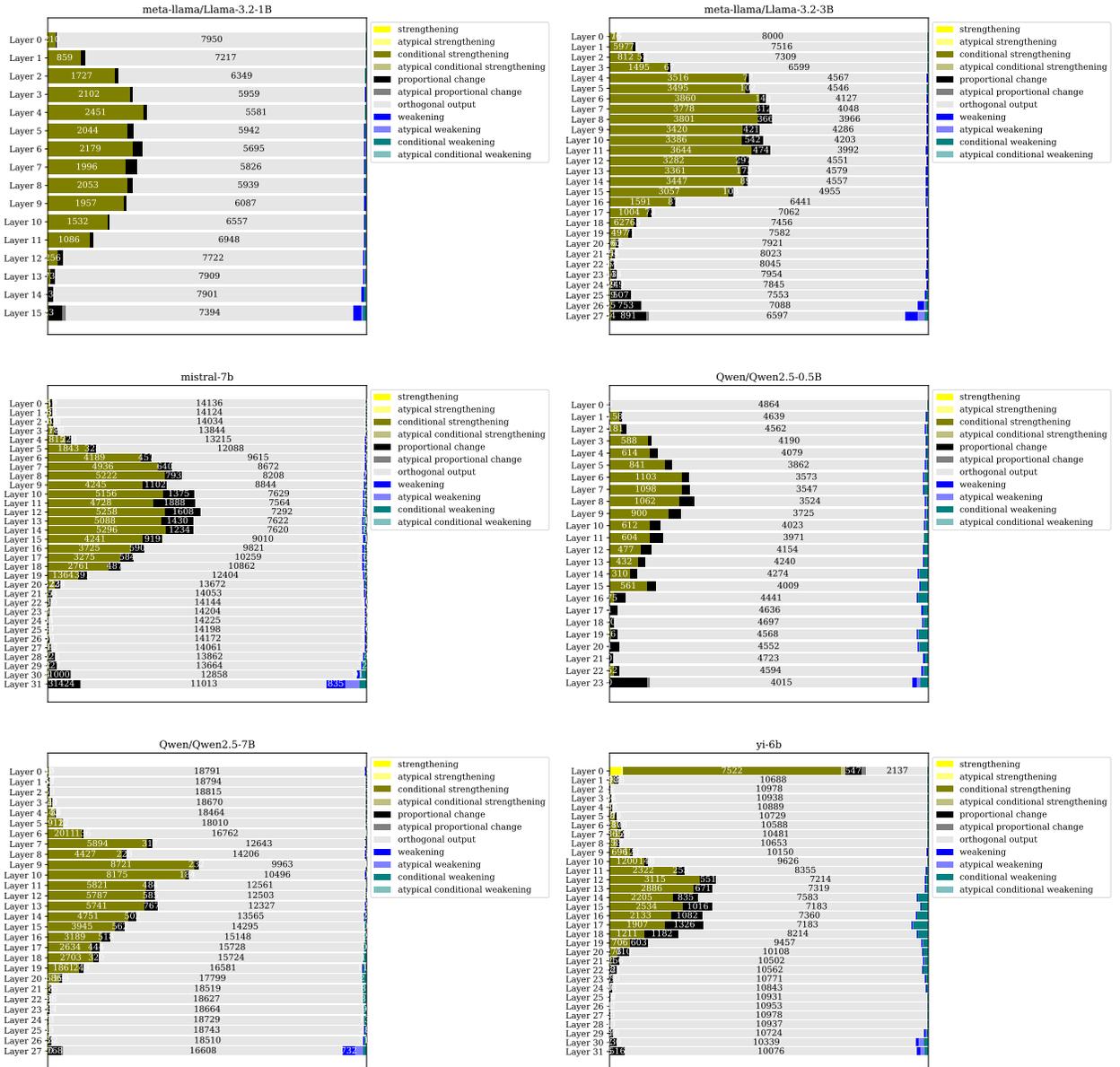


Figure 41: Distribution of neurons by layer and category for a range of models

2538  
2539  
2540  
2541  
2542  
2543  
2544  
2545  
2546  
2547  
2548  
2549  
2550  
2551  
2552  
2553  
2554  
2555  
2556  
2557  
2558  
2559  
2560  
2561  
2562  
2563  
2564  
2565  
2566  
2567  
2568  
2569  
2570  
2571  
2572  
2573  
2574  
2575  
2576  
2577  
2578  
2579  
2580  
2581  
2582  
2583  
2584  
2585  
2586  
2587  
2588  
2589  
2590  
2591



2592  
 2593  
 2594  
 2595  
 2596  
 2597  
 2598  
 2599  
 2600  
 2601  
 2602  
 2603  
 2604  
 2605  
 2606  
 2607  
 2608  
 2609  
 2610  
 2611  
 2612  
 2613  
 2614  
 2615  
 2616  
 2617  
 2618  
 2619  
 2620  
 2621  
 2622  
 2623  
 2624  
 2625  
 2626  
 2627  
 2628  
 2629  
 2630  
 2631  
 2632  
 2633  
 2634  
 2635  
 2636  
 2637  
 2638  
 2639  
 2640  
 2641  
 2642  
 2643  
 2644  
 2645

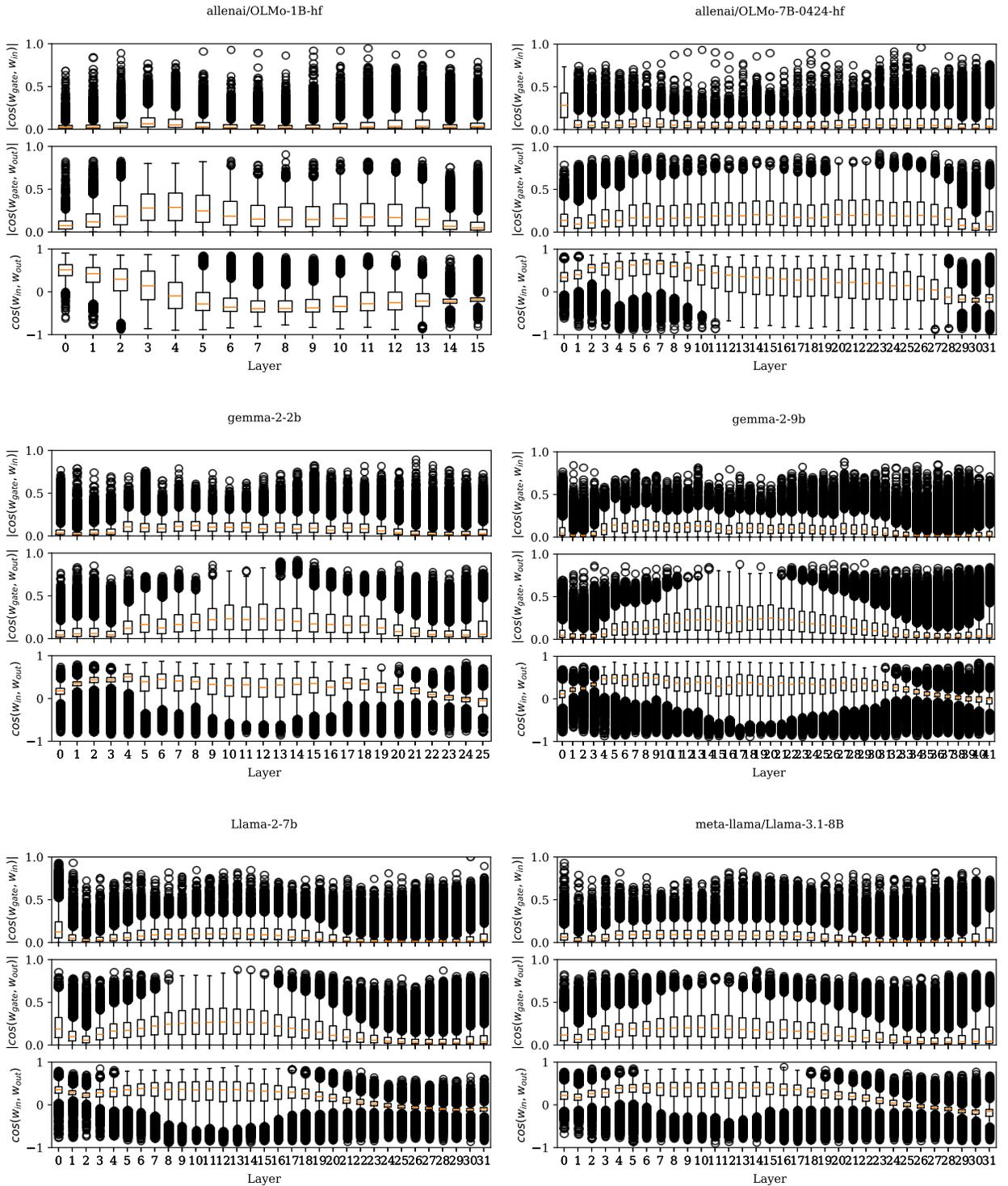


Figure 43: Boxplots for the distribution of weight cosine similarities in each layer. For  $\cos(w_{gate}, w_{in})$  and  $\cos(w_{gate}, w_{out})$  we show the absolute value since their sign does not carry any information on its own.

2646  
2647  
2648  
2649  
2650  
2651  
2652  
2653  
2654  
2655  
2656  
2657  
2658  
2659  
2660  
2661  
2662  
2663  
2664  
2665  
2666  
2667  
2668  
2669  
2670  
2671  
2672  
2673  
2674  
2675  
2676  
2677  
2678  
2679  
2680  
2681  
2682  
2683  
2684  
2685  
2686  
2687  
2688  
2689  
2690  
2691  
2692  
2693  
2694  
2695  
2696  
2697  
2698  
2699

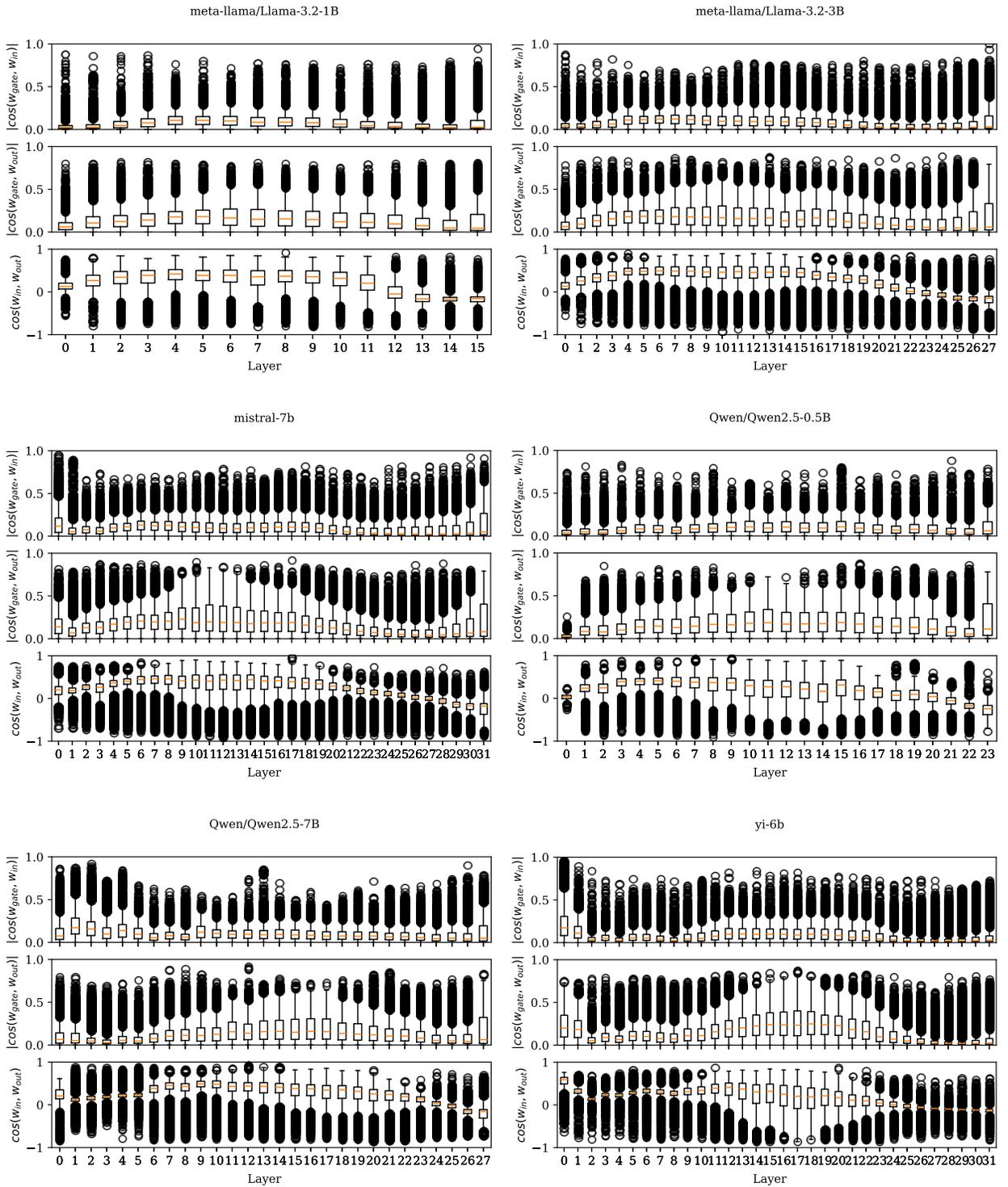


Figure 44: Continuation of figure 43.

2700  
2701  
2702  
2703  
2704  
2705  
2706  
2707  
2708  
2709  
2710  
2711  
2712  
2713  
2714  
2715  
2716  
2717  
2718  
2719  
2720  
2721  
2722  
2723  
2724  
2725  
2726  
2727  
2728  
2729  
2730  
2731  
2732  
2733  
2734  
2735  
2736  
2737  
2738  
2739  
2740  
2741  
2742  
2743  
2744  
2745  
2746  
2747  
2748  
2749  
2750  
2751  
2752  
2753

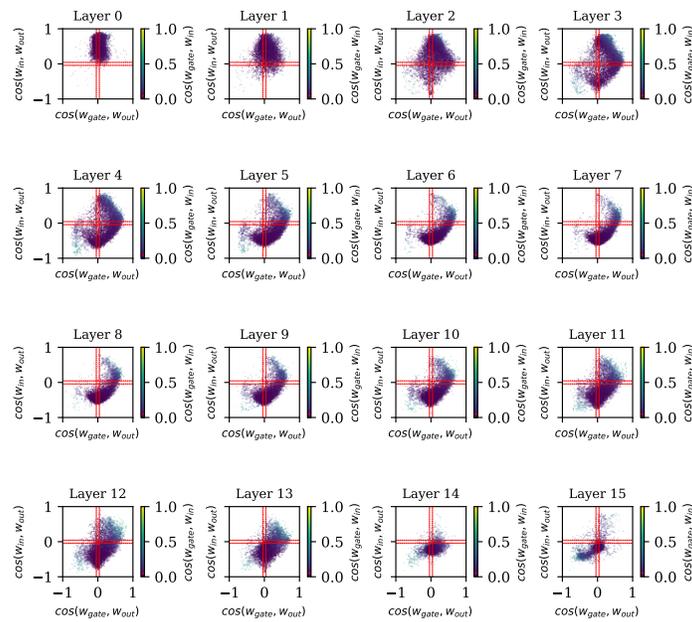


Figure 45: Equivalent of figure 2 for OLMo-1B

2754  
2755  
2756  
2757  
2758  
2759  
2760  
2761  
2762  
2763  
2764  
2765  
2766  
2767  
2768  
2769  
2770  
2771  
2772  
2773  
2774  
2775  
2776  
2777  
2778  
2779  
2780  
2781  
2782  
2783  
2784  
2785  
2786  
2787  
2788  
2789  
2790  
2791  
2792  
2793  
2794  
2795  
2796  
2797  
2798  
2799  
2800  
2801  
2802  
2803  
2804  
2805  
2806  
2807

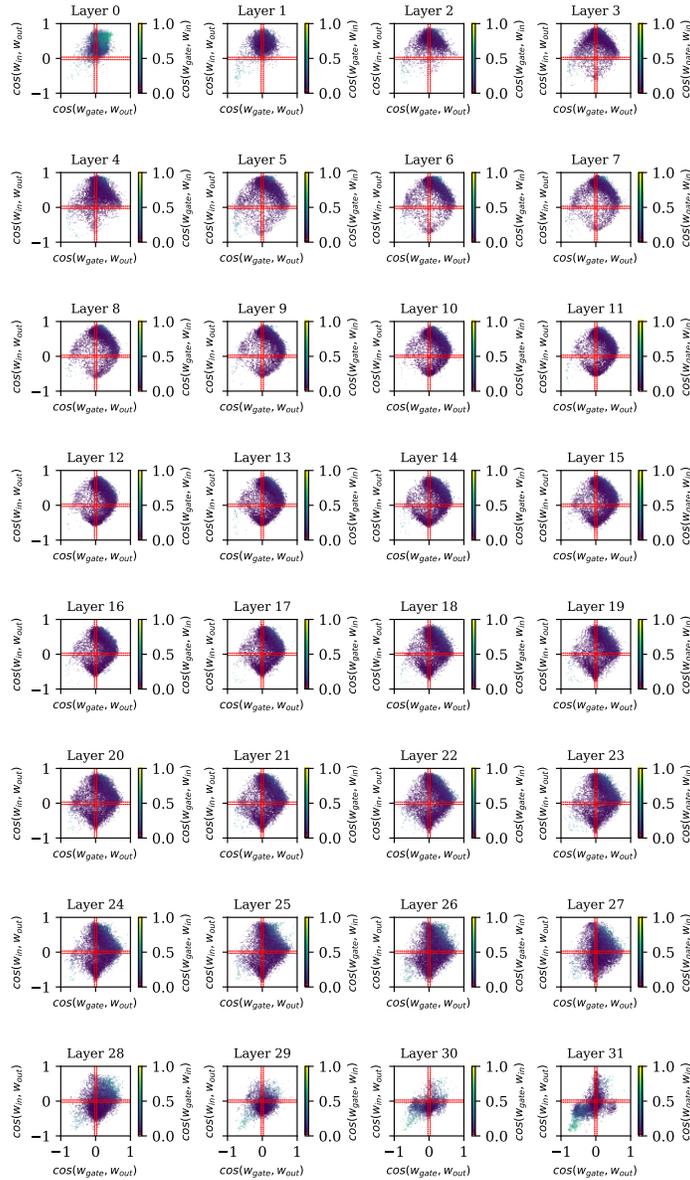


Figure 46: Equivalent of figure 2 for OLMo-7B-0424

2808  
 2809  
 2810  
 2811  
 2812  
 2813  
 2814  
 2815  
 2816  
 2817  
 2818  
 2819  
 2820  
 2821  
 2822  
 2823  
 2824  
 2825  
 2826  
 2827  
 2828  
 2829  
 2830  
 2831  
 2832  
 2833  
 2834  
 2835  
 2836  
 2837  
 2838  
 2839  
 2840  
 2841  
 2842  
 2843  
 2844  
 2845  
 2846  
 2847  
 2848  
 2849  
 2850  
 2851  
 2852  
 2853  
 2854  
 2855  
 2856  
 2857  
 2858  
 2859  
 2860  
 2861

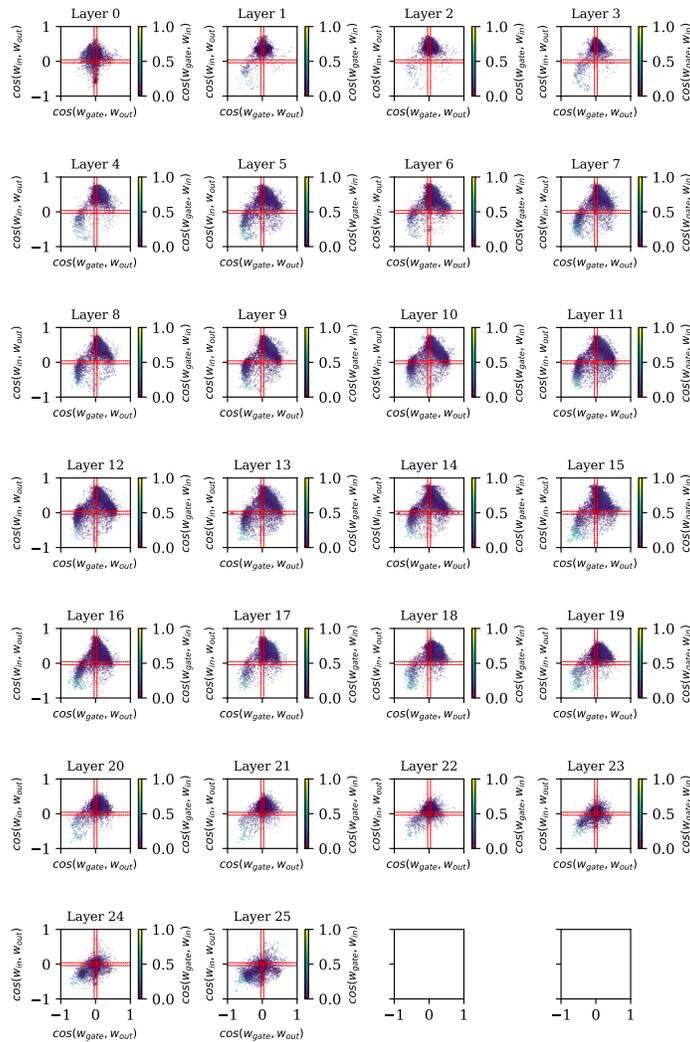


Figure 47: Equivalent of figure 2 for Gemma-2-2B

2862  
2863  
2864  
2865  
2866  
2867  
2868  
2869  
2870  
2871  
2872  
2873  
2874  
2875  
2876  
2877  
2878  
2879  
2880  
2881  
2882  
2883  
2884  
2885  
2886  
2887  
2888  
2889  
2890  
2891  
2892  
2893  
2894  
2895  
2896  
2897  
2898  
2899  
2900  
2901  
2902  
2903  
2904  
2905  
2906  
2907  
2908  
2909  
2910  
2911  
2912  
2913  
2914  
2915

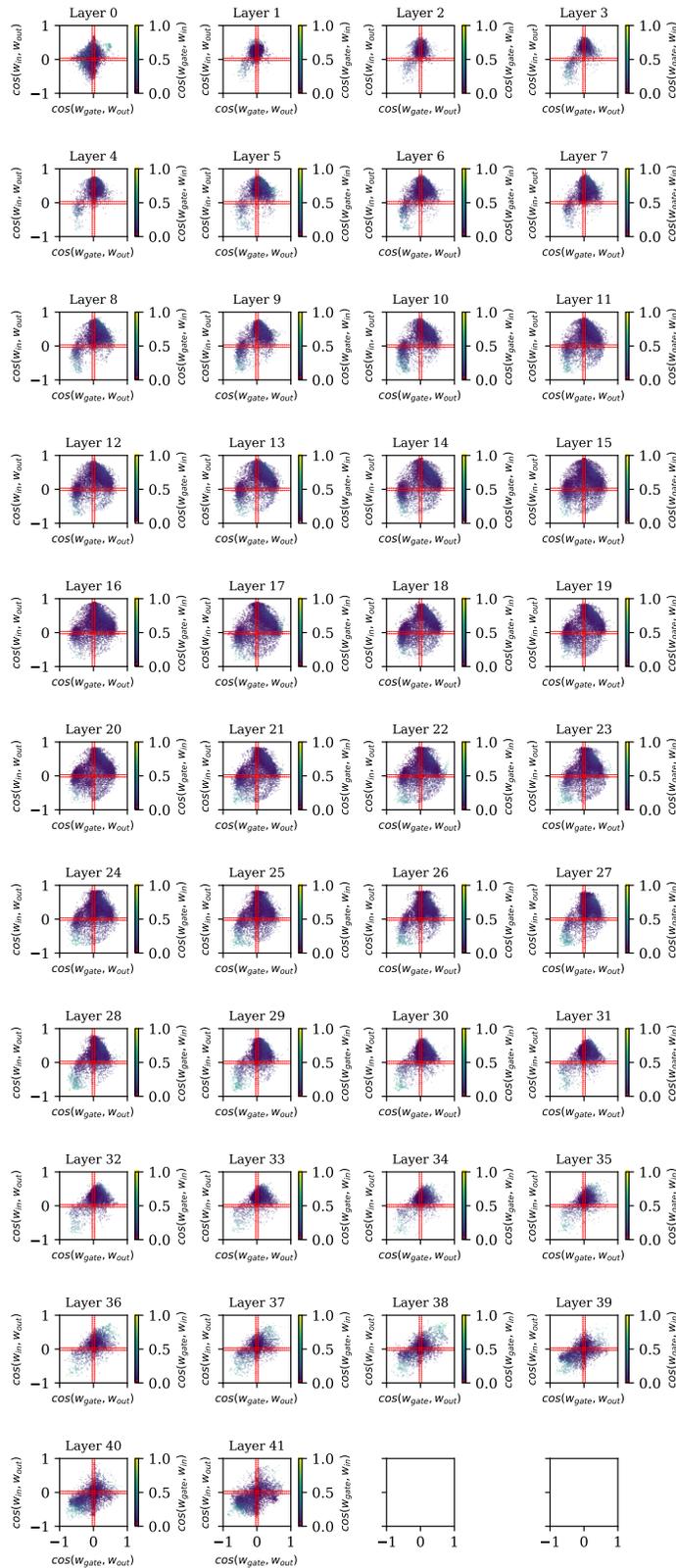


Figure 48: Equivalent of figure 2 for Gemma-2-9B

2916  
2917  
2918  
2919  
2920  
2921  
2922  
2923  
2924  
2925  
2926  
2927  
2928  
2929  
2930  
2931  
2932  
2933  
2934  
2935  
2936  
2937  
2938  
2939  
2940  
2941  
2942  
2943  
2944  
2945  
2946  
2947  
2948  
2949  
2950  
2951  
2952  
2953  
2954  
2955  
2956  
2957  
2958  
2959  
2960  
2961  
2962  
2963  
2964  
2965  
2966  
2967  
2968  
2969

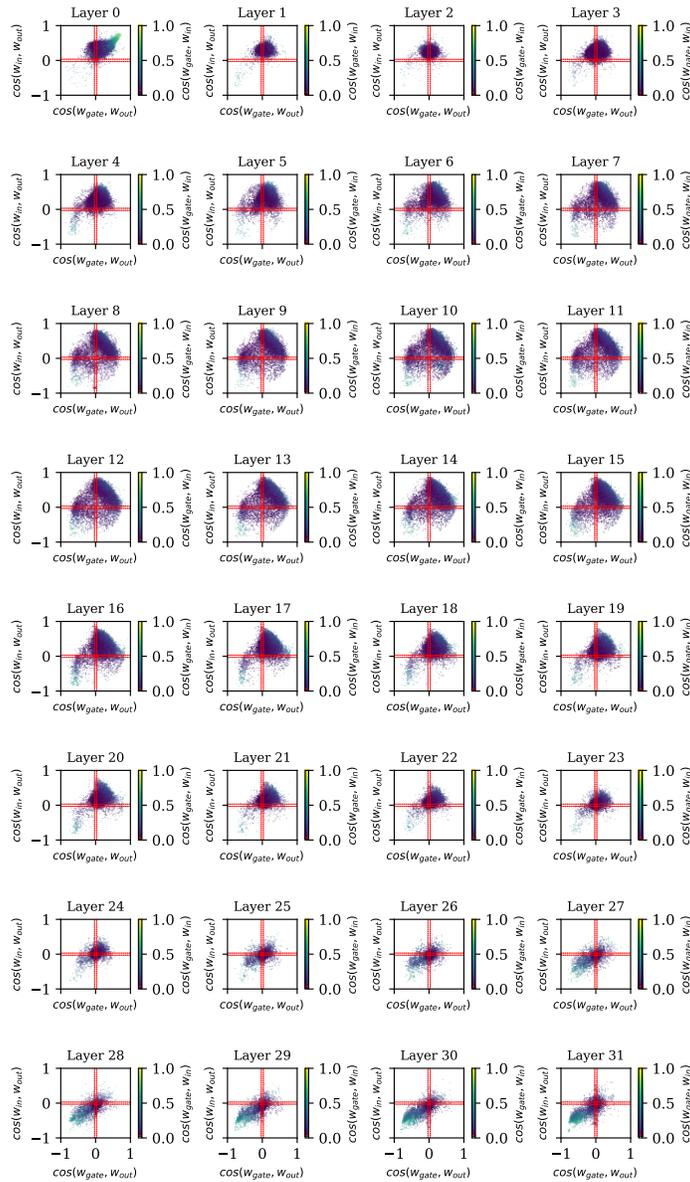


Figure 49: Equivalent of figure 2 for Llama-2-7B

2970  
 2971  
 2972  
 2973  
 2974  
 2975  
 2976  
 2977  
 2978  
 2979  
 2980  
 2981  
 2982  
 2983  
 2984  
 2985  
 2986  
 2987  
 2988  
 2989  
 2990  
 2991  
 2992  
 2993  
 2994  
 2995  
 2996  
 2997  
 2998  
 2999  
 3000  
 3001  
 3002  
 3003  
 3004  
 3005  
 3006  
 3007  
 3008  
 3009  
 3010  
 3011  
 3012  
 3013  
 3014  
 3015  
 3016  
 3017  
 3018  
 3019  
 3020  
 3021  
 3022  
 3023

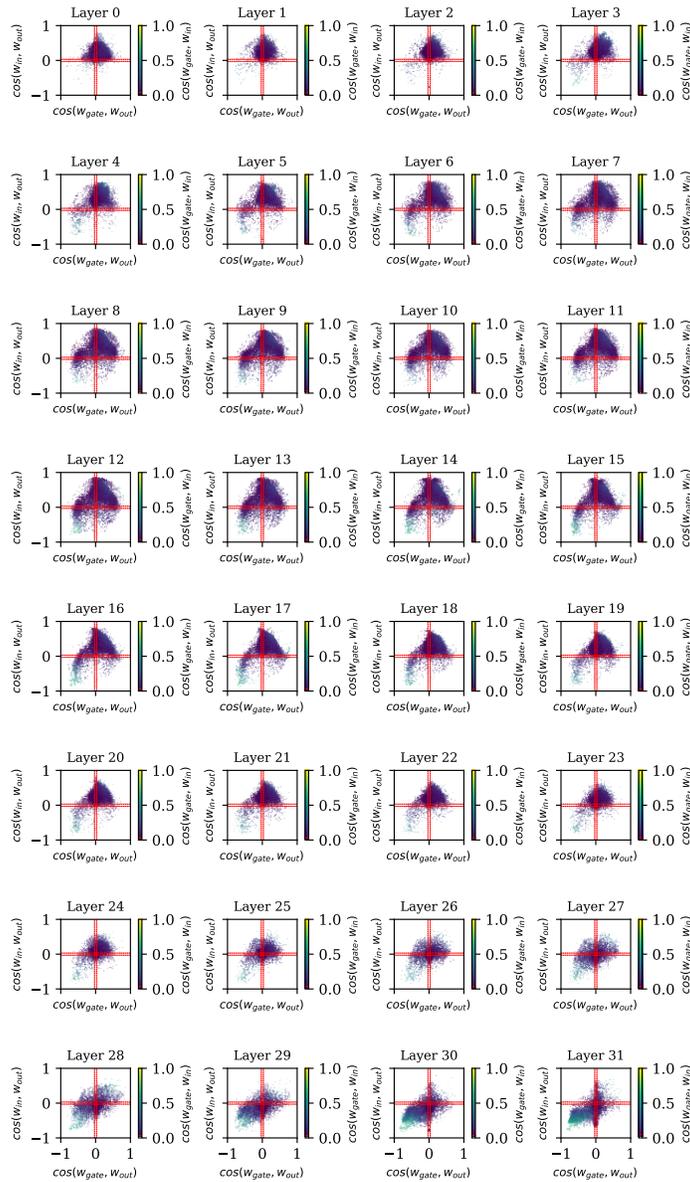


Figure 50: Equivalent of figure 2 for Llama-3.1-8B

3024  
3025  
3026  
3027  
3028  
3029  
3030  
3031  
3032  
3033  
3034  
3035  
3036  
3037  
3038  
3039  
3040  
3041  
3042  
3043  
3044  
3045  
3046  
3047  
3048  
3049  
3050  
3051  
3052  
3053  
3054  
3055  
3056  
3057  
3058  
3059  
3060  
3061  
3062  
3063  
3064  
3065  
3066  
3067  
3068  
3069  
3070  
3071  
3072  
3073  
3074  
3075  
3076  
3077

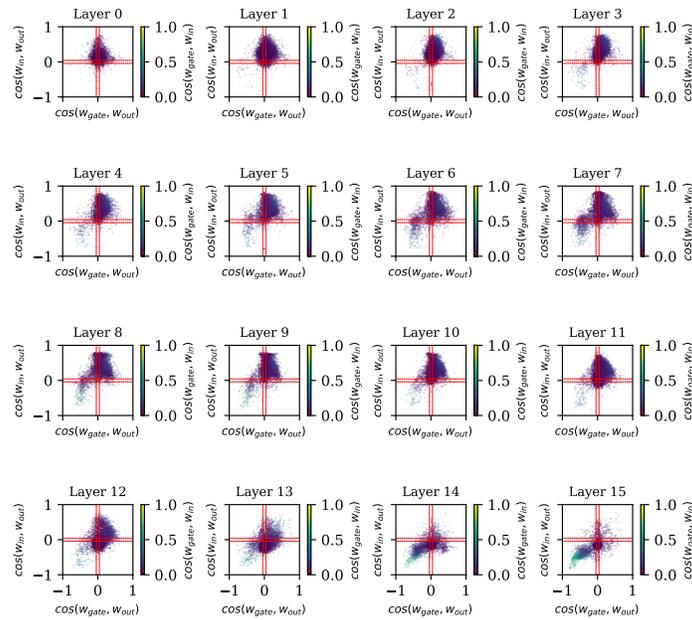


Figure 51: Equivalent of figure 2 for Llama-3.2-1B

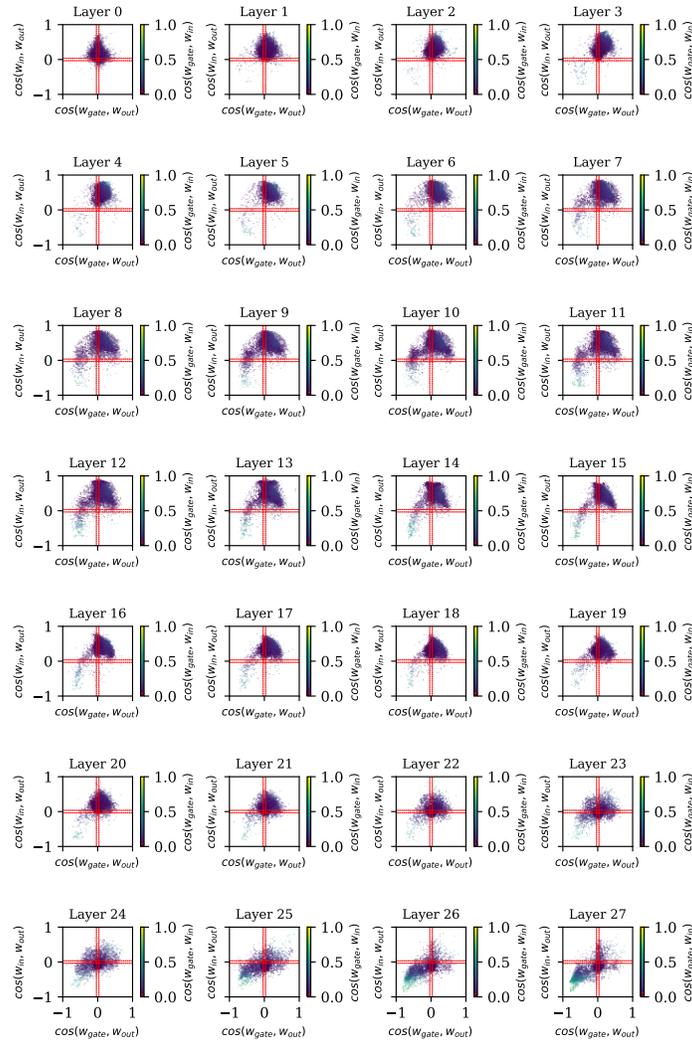


Figure 52: Equivalent of figure 2 for Llama-3.2-3B (same model but all layers)

3132  
3133  
3134  
3135  
3136  
3137  
3138  
3139  
3140  
3141  
3142  
3143  
3144  
3145  
3146  
3147  
3148  
3149  
3150  
3151  
3152  
3153  
3154  
3155  
3156  
3157  
3158  
3159  
3160  
3161  
3162  
3163  
3164  
3165  
3166  
3167  
3168  
3169  
3170  
3171  
3172  
3173  
3174  
3175  
3176  
3177  
3178  
3179  
3180  
3181  
3182  
3183  
3184  
3185

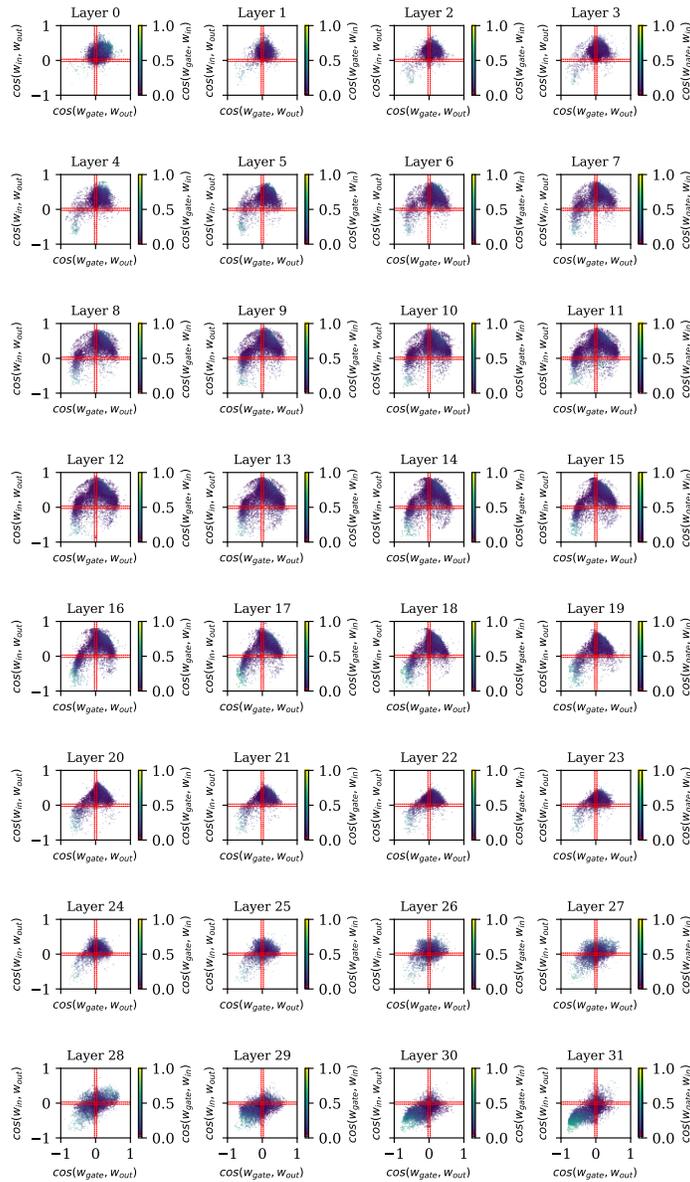


Figure 53: Equivalent of figure 2 for Mistral-7B

3186  
3187  
3188  
3189  
3190  
3191  
3192  
3193  
3194  
3195  
3196  
3197  
3198  
3199  
3200  
3201  
3202  
3203  
3204  
3205  
3206  
3207  
3208  
3209  
3210  
3211  
3212  
3213  
3214  
3215  
3216  
3217  
3218  
3219  
3220  
3221  
3222  
3223  
3224  
3225  
3226  
3227  
3228  
3229  
3230  
3231  
3232  
3233  
3234  
3235  
3236  
3237  
3238  
3239

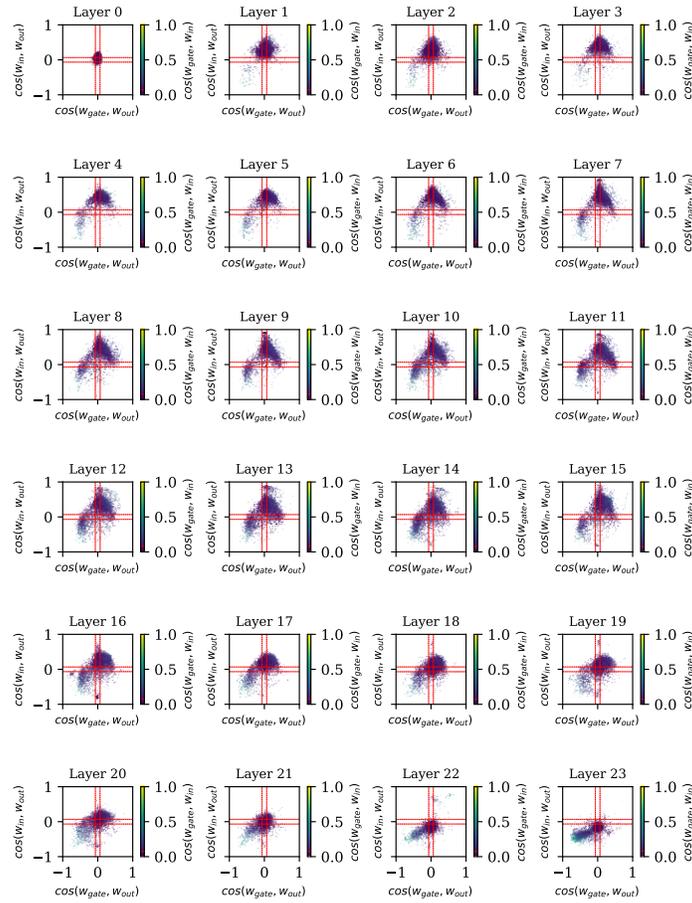


Figure 54: Equivalent of figure 2 for Qwen2.5-0.5B

3240  
3241  
3242  
3243  
3244  
3245  
3246  
3247  
3248  
3249  
3250  
3251  
3252  
3253  
3254  
3255  
3256  
3257  
3258  
3259  
3260  
3261  
3262  
3263  
3264  
3265  
3266  
3267  
3268  
3269  
3270  
3271  
3272  
3273  
3274  
3275  
3276  
3277  
3278  
3279  
3280  
3281  
3282  
3283  
3284  
3285  
3286  
3287  
3288  
3289  
3290  
3291  
3292  
3293

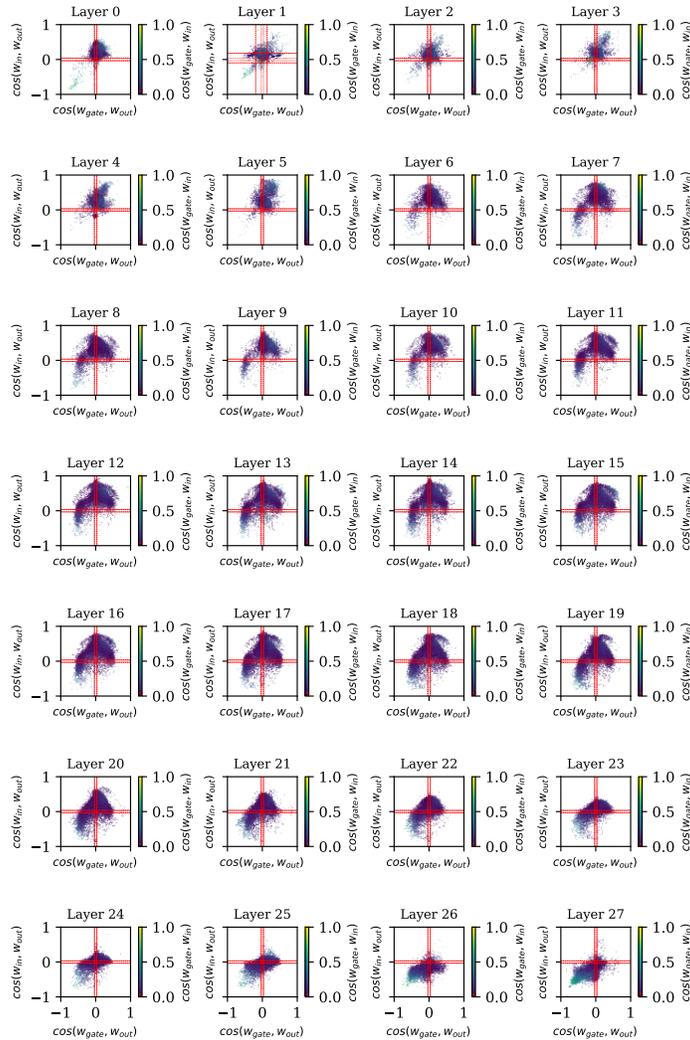


Figure 55: Equivalent of figure 2 for Qwen2.5-7B

3294  
 3295  
 3296  
 3297  
 3298  
 3299  
 3300  
 3301  
 3302  
 3303  
 3304  
 3305  
 3306  
 3307  
 3308  
 3309  
 3310  
 3311  
 3312  
 3313  
 3314  
 3315  
 3316  
 3317  
 3318  
 3319  
 3320  
 3321  
 3322  
 3323  
 3324  
 3325  
 3326  
 3327  
 3328  
 3329  
 3330  
 3331  
 3332  
 3333  
 3334  
 3335  
 3336  
 3337  
 3338  
 3339  
 3340  
 3341  
 3342  
 3343  
 3344  
 3345  
 3346  
 3347

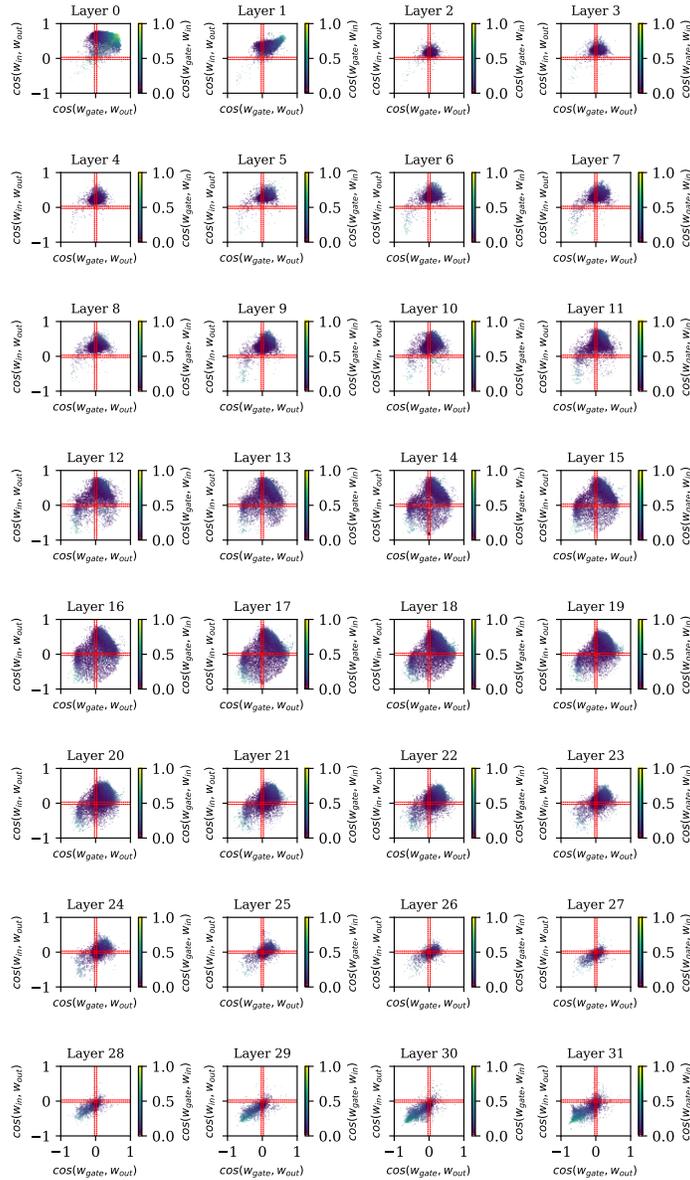


Figure 56: Equivalent of figure 2 for Yi-6B