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 an analysis of their read-write functionality, 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 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. ¹

1 Introduction and background

Despite recent progress in interpretability, there is still much that is unclear about how transformer-based (Vaswani et al., 2017) large language models (LLMs) achieve their impressive performance. Prior interpretability work has addressed the interpretation of MLP sublayers, and we follow this line of research. We expand this work in two directions: First, much previous work analyzes neurons based only on the contexts in which they activate (Voita et al., 2024) or based only on their output weights² (Gurnee et al., 2024). In contrast (inspired by ideas from Elhage et al. (2021); Gurnee et al. (2024)), we put the read-write (RW) functionality of neurons in the center of our analysis, and classify neurons according to the interactions between input and output weights. Second, we do this in the context of 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.

This new approach allows us to gain a number of novel insights into the inner workings of these models. 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) 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.

¹We publish code at https://anonymous.4open.science/r/RW_functionalities-4D32.

²We use "weight" to refer to a weight vector, not a scalar.

062 063 064

065 067

068

069

071

074 075 076

077 078

073

079 080 081

084 085

082

090

092

094 095 096

098

099

101 102

103 104 105

106 107

Table 1: Six prototypical read-write (RW) functionalities.

	$ \cos(oldsymbol{w}_{ ext{gate}}, oldsymbol{w}_{ ext{out}}) $	$\gg 0$	≈ 0
$\cos(oldsymbol{w}_{ m in}, oldsymbol{w}_{ m out})$			
$\gg 0$		strengthening	conditional
			strengthening
$\ll 0$		weakening	conditional
			weakening
≈ 0		proportional	orthogonal
		change	output

THEORETICAL FRAMEWORK

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, it detects certain directions in the residual stream (i.e., whether the current residual stream vector at least approximately points in one of these directions in model space), corresponding to its weight vectors on the input side; and then writes to a certain direction, corresponding to its output weight vector.

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.

1.2 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. We describe single neurons as opposed to whole MLP layers.

Traditional activation functions like ReLU take a single scalar as argument: $x_{post} = ReLU(x_{in})$. In contrast, a gated activation function like SwiGLU takes two arguments:

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

The scalars x_{gate} and x_{in} are the dot products of the (normalized) residual stream (which we denote x_{norm}) with the neuron's **gate** and **linear input** weight vectors, w_{gate} and w_{in} .

Finally, $x_{\text{post}} w_{\text{out}}$ (a rescaled version of the **output** weight vector) is added to the residual stream.

Fully written out: The neuron adds

$$Swish(\boldsymbol{x}_{norm} \cdot \boldsymbol{w}_{gate}) \cdot (\boldsymbol{x}_{norm} \cdot \boldsymbol{w}_{in}) \cdot \boldsymbol{w}_{out} \tag{1}$$

to the residual stream.

METHOD

WEIGHT PREPROCESSING

We note two closely related theoretical properties of gated activation functions:

- 1. Positive vs negative activation. Strong activations can be either positive or negative: If $x_{\text{gate}} \gg 0$ (i.e. $x_{\text{norm}} \cdot w_{\text{gate}} \gg 0$), then x_{in} (i.e. $x_{\text{norm}} \cdot w_{\text{in}}$) can still be strongly positive or negative, because $w_{\rm in}$ is different from $w_{\rm gate}$. In the case $x_{\rm gate} \gg 0$, $x_{\rm in} \gg 0$, we get a strong positive activation. In the case $x_{\text{gate}} \gg 0$, $x_{\text{in}} \ll 0$, we get a strong negative activation. So, depending on the context, a given gated activation neuron can either add the output weight vector to the residual stream or subtract it.
- **2. Symmetry.** Switching the signs of both w_{in} and w_{out} preserves behavior: In equation (1), if we replace $w_{\rm in}$ by $-w_{\rm in}$ and $w_{\rm out}$ by $-w_{\rm out}$, the two minus signs cancel out.

We propose a **weight preprocessing** step to make interpretation of gated neurons more intuitive. Specifically, we multiply w_{in} and w_{out} by the sign of $\cos(w_{gate}, w_{in})$. By property (2) above, this does not change neuron behavior.

The effect of this operation is that 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. Additionally, it guarantees that two equivalent neurons by property (2) will now also have the same weights.

2.2 READ-WRITE FUNCTIONALITIES

Our main research question is: What is the relationship between what a neuron reads and what it writes? We address this question by computing the cosine similarity of input and output weights. We present a taxonomy of prototypical RW functionalities in table 1.

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: $\cos(w_{\rm in}, w_{\rm out}) \approx \pm 1$. The neuron detects a direction and then adds it to / removes it from the residual stream. (ii) *Conditional strengthening* / *weakening* neurons: $w_{\rm in}$ and $w_{\rm out}$ are roughly collinear and $w_{\rm gate}$ is orthogonal to them. The neuron also strengthens / weakens the direction detected by its $w_{\rm in}$ vector, but will only activate *conditional on* $w_{\rm gate}$ being present in the residual stream. (iii) *Proportional change* neurons: $w_{\rm out}$ is collinear to $w_{\rm gate}$, but is orthogonal to $w_{\rm in}$. If $w_{\rm gate}$ 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 $w_{\rm in}$ in the residual stream.

These prototypical classes are limited in scope: Many cosines will not be close to 0 or ± 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 (1), $\cos(\boldsymbol{w}_{\text{in}}, \boldsymbol{w}_{\text{gate}})$ may not always "match" the other two cosine similarities. For example, the two reading weights may be orthogonal, but $\boldsymbol{w}_{\text{out}} = \frac{1}{\sqrt{2}} \boldsymbol{w}_{\text{gate}} + \frac{1}{\sqrt{2}} \boldsymbol{w}_{\text{in}}$; then both cosine similarities are $\frac{1}{\sqrt{2}} > 0.5$. We are mainly interested in RW behavior rather than comparing the two reading weights, so we classify such cases based on $\cos(\boldsymbol{w}_{\text{in}}, \boldsymbol{w}_{\text{out}})$ and $\cos(\boldsymbol{w}_{\text{gate}}, \boldsymbol{w}_{\text{out}})$. To signal the "mismatch" of $\cos(\boldsymbol{w}_{\text{in}}, \boldsymbol{w}_{\text{gate}})$, we prepend *atypical* to the category's name. In the above example, we will speak of an atypical strengthening neuron.

2.3 RANDOM BASELINES

Given a cosine similarity between weight vectors, we want to test if it is significantly different from random. To do so, we consider two different random baselines: (i) i.i.d. Gaussian vector entries (as in a randomly initialized model), or (ii) a layer-specific baseline based on "mismatched cosines". We describe both approaches in appendix C.

In practice, we find that both baselines give quite similar 95% randomness ranges.

3 Where to find weakening neurons

In this section we 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.

Concretely, we apply our method to 12 LLMs: Gemma-2-2B, Gemma-2-9B Gemma (2024), Llama-2-7B, Llama-3.1-8B, Llama-3.2-1B, Llama-3.2-3B Touvron et al. (2023), OLMo-1B, OLMo-7B-0424 (Groeneveld et al., 2024), Mistral-7B (Jiang et al., 2023), Qwen2.5-0.5B, Qwen2.5-7B (Yang et al.,

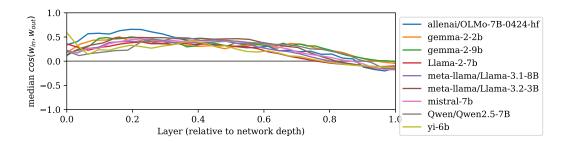


Figure 1: Median of $\cos(w_{\rm in}, w_{\rm out})$ by layer (x-axis) for 9 models of 2B to 9B parameters. For all models, the value is positive in the beginning and negative in the end, indicating that early-middle layers "strengthen" directions they find in the residual stream whereas later layers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

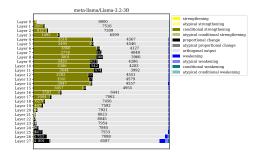


Figure 2: Distribution of neurons by layer and category.

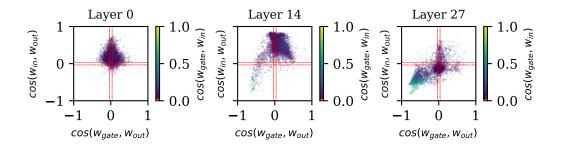


Figure 3: Fine-grained analysis of neuron RW behavior in three layers of Llama-3.2-3B, based on the configuration of their three weight vectors in parameter space. Each subplot represents a layer, each dot a neuron. The red lines mark the 95% randomness regions for each of the three cosine values. (There is a dotted line for variant (i) and a dashed line for variant (ii) in section 2.3, 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.

2024), Yi-6B (01.AI et al., 2025). These models use SwiGLU, except for Gemma, which uses GeGLU.

To demonstrate our finding in more detail, we present three representative plots. (See appendix G for more.) Figure 1 shows the median value of $\cos(\boldsymbol{w}_{\text{in}}, \boldsymbol{w}_{\text{out}})$ across all layers of the nine larger models (2B and above). Here the common pattern is especially visible: In early-middle layers of all models, a majority of neurons has a $\cos(\boldsymbol{w}_{\text{in}}, \boldsymbol{w}_{\text{out}})$ high above zero, indicating strengthening; in late layers, still for all models, this median cosine similarity goes slightly below zero, indicating a majority of weakening neurons.

Our other two plots focus on **Llama-3.2-3B**, but appendix G contains many similar plots for other models, and the patterns we describe here are universal. Figure 2 shows RW class distribution across layers. In figure 3, we plot the distribution of neurons in a few selected layers, by displaying each neuron as a point with $\cos(\boldsymbol{w}_{\text{gate}}, \boldsymbol{w}_{\text{out}})$ indicated on the x-axis, $\cos(\boldsymbol{w}_{\text{in}}, \boldsymbol{w}_{\text{out}})$ on the y-axis and $\cos(\boldsymbol{w}_{\text{gate}}, \boldsymbol{w}_{\text{in}})$ as its color.

Input manipulation. First, we see that a large proportion of neurons are input manipulators (i.e., they are not orthogonal output neurons): In figure 2, these are 25% of all neurons, and as much as 50% in early-middle layers (layers $7-11^3$). What is more, figure 3 shows that even the neurons classified as "orthogonal output" often belong to clusters that are centered above/below the horizontal line. Their weight cosine similarities often exceed the significance threshold. E.g., in layer 14, there are many neurons whose $\cos(w_{\rm in}, w_{\rm out})$ (y-axis) is below 0.5 but above the significance threshold. This suggests that even the "orthogonal output" neurons perform input manipulation to some extent. This finding is universal.

Different RW functionalities. Weakening neurons represent a large share of the (relatively few) input manipulators in late layers. They form a somewhat separate cluster in figure 3 (in the bottom-left corner of the rightmost subplot). Another important input manipulator class in late layers is proportional change. In contrast, across all models, early middle layers are dominated by conditional strengthening. In fact, the majority of input manipulators (more than 80% in Llama) belong to just this one class.

This general pattern of strengthening-then-weakening holds across models, as figure 1 shows at one glance. In figure 3 (and figure 34 in the appendix), the pattern manifests as a large cluster of neurons, centered clearly above the x-axis in most layers, but moving below it in the last few layers.

While the above patterns are universal, some other models display additional patterns: for example OlMo-7B and especially OLMo-1B display a large number of conditional weakening neurons in middle-late layers. See appendix G.

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

4 ABLATION EXPERIMENTS

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

In this and the following sections, the experiments involve running a model on a dataset. Therefore, to save resources, we focus on a single model: We choose **OLMo-7B**, because its training dataset, Dolma (Soldaini et al., 2024), is publicly available and its RW functionalities mostly follow the typical patterns. As a dataset, we use a random subset of 20M tokens from Dolma,⁴ except for attribute rate, where we follow the setup of Geva et al. (2023).

³We use zero-based (Python-style) indexing throughout the paper.

⁴The size of 20M tokens follows Voita et al., 2024.

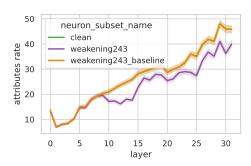


Figure 4: Effect on attribute rate of ablating all 243 weakening neurons (weakening243), or 243 random neurons from the same layers (weakening243_baseline). The effect of weakening neurons is clearly visible, already from layer ≈ 10 onward, even though weakening neurons are few and mostly in late layers.

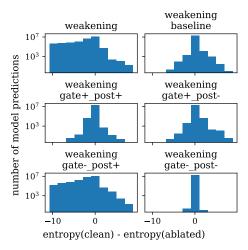


Figure 5: Effect on entropy of various neuron activations. E.g., in $\approx 10^5$ next-token predictions, weakening neurons decrease the entropy by about 10 nats, whereas they increase it much more rarely. "Weakening baseline" denotes random neurons from the same layers as weakening neurons. The bottom four plots describe the results of conditional ablations (section 4.2). E.g., "gate+_post+" describes the effect of those activations in which $x_{\rm gate} > 0$ and $x_{\rm post} > 0$.

4.1 EXPLORING RW CLASSES AND METRICS

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

The main metrics we consider are attribute rate (Geva et al., 2023) and entropy of the output distribution. For simplicity we choose zero ablation, i.e., ablating a neuron means setting its activation to zero. We justify our choices in appendix D.

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.

For attribute rate, the effect is most visible in layers ≈ 10 and onward. See figure 4. This is particularly interesting since there are very few weakening neurons in these early-middle layers.

The case of entropy is also striking: Figure 5 shows that ablating weakening neurons 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.

4.2 CONDITIONAL ABLATIONS

We now further investigate the effect of weakening neurons on entropy. We use **conditional ablations**: We zero out only some activations of each neuron, based on the signs of the corresponding $x_{\rm gate}$ and $x_{\rm in}$. Specifically, we consider the following four conditions (the definitions assume our weight preprocessing, cf. section 2.1): (i) $x_{\rm gate} > 0$, $x_{\rm in} > 0$, leading to $x_{\rm post} > 0$; (ii) $x_{\rm gate} > 0$, $x_{\rm in} < 0$, leading to $x_{\rm post} < 0$; (iii) $x_{\rm gate} < 0$, $x_{\rm in} < 0$, leading to $x_{\rm post} < 0$; (iv) $x_{\rm gate} < 0$, $x_{\rm in} > 0$, leading to $x_{\rm post} < 0$.

We find that the sharpening effect of weakening neurons is largely due to case (iii): In figure 5, only case (iii) (bottom left subplot) shows entropy effects similar to those of weakening neurons as a whole. 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 sharpen it).

It is *surprising* for two reasons: First, these negative $x_{\rm gate}$ activations are relatively rare in weakening neurons (as we will see in section 5). Second, **negative gate values** are relatively small (whereas positive values can be arbitrarily large). Researchers previously hypothesized that these negative activations might play a role in the internal mechanisms of neural networks (see e.g. Gurnee et al., 2023^5), but so far this has been speculation only. Our results show for the first time that this effect indeed occurs empirically and that it can have a strong effect on the behavior of transformer-based models.

Our finding 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. (In the prototypical case $w_{\text{out}} \approx -w_{\text{in}}$, when the neuron detects the presence of w_{in} , it will now write $-w_{\text{out}} = -w_{\text{in}} = w_{\text{in}}$ instead of w_{in} .) Thus, in these cases, a neuron can further increase the score of a high-scoring token, or decrease the score of a token with large negative score, and thus indeed make the output distribution sharper.

5 WEAKENING NEURONS ACTIVATE OFTEN

Our findings from section 4 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 $\cos(\boldsymbol{w}_{in}, \boldsymbol{w}_{out})$ – but in a GELU model. We now investigate whether a similar phenomenon occurs with gated activation functions.

⁵They assumed that GELU (or equivalently, Swish) is "essentially the same activation as a ReLU", and said they "would be particularly excited to see future work exhibiting [...] case studies" of mechanisms involving negative values of such an activation function.

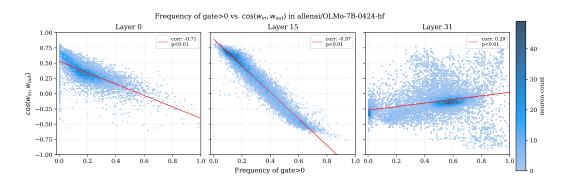


Figure 6: Two-dimensional histogram of activation frequency (x-axis) vs. $\cos(w_{\text{in}}, w_{\text{out}})$ (y-axis).

For the results, see figure 6 and additional figures and tables in appendix 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 $\cos(\mathbf{w}_{\text{in}}, \mathbf{w}_{\text{out}})$ and activation frequency: correlations are at least -0.71 in all layers except the last two (which have -0.29 and +0.29).

The last layer displays a different pattern (rightmost subplot in figure 6). Here the correlation is positive (+0.29), and we can distinguish two clusters of neurons: One cluster has a medium-negative $\cos(\boldsymbol{w}_{\text{in}}, \boldsymbol{w}_{\text{out}})$ (around -0.3) and activates very rarely; another one is much more spread out (both in terms of $\cos(\boldsymbol{w}_{\text{in}}, \boldsymbol{w}_{\text{out}})$ 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.

It seems that each conditional strengthening neuron is responsible only for a narrow domain, perhaps a specific set of tokens. In any case, 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 4).

6 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) compute a singular value decomposition (SVD) of layer weights and often 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. Taking the middle ground, Gurnee et al. (2023) argue that interpretable features correspond to sparse combinations of neurons, but this includes 1-sparse combinations, i.e., individual neurons.

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. They note that suppression neurons seem to activate "when it is plausible but not certain that the next token is from the relevant set". Stolfo et al. (2024) also investigate some output-based neuron classes.

Researchers have paid less attention to 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 (negative cosines between input and output weights "may also be mechanisms for conditionally deleting information", footnote 7), but do not follow up on this remark. Note also 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.

7 CONCLUSION

We explore the read-write perspective for investigating gated neurons in LLMs. Our method complements prior interpretability approaches and provides new insights into the inner workings of LLMs.

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

Beyond weakening neurons, we have observed that a large share of neurons exhibit other significant RW interactions. In particular, early-middle layers are dominated by conditional strengthening neurons. This finding is particularly significant since it is universal across models.

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

REFERENCES

01.AI, :, Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Guoyin Wang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yanpeng Li, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. Yi: Open foundation models by 01.ai, 2025. URL https://arxiv.org/abs/2403.04652.

Nora Belrose, Zach Furman, Logan Smith, Danny Halawi, Igor Ostrovsky, Lev McKinney, Stella Biderman, and Jacob Steinhardt. Eliciting latent predictions from transformers with the tuned lens, 2023. URL https://arxiv.org/pdf/2303.08112.

Fahim Dalvi, Nadir Durrani, Hassan Sajjad, Yonatan Belinkov, Anthony Bau, and James Glass. What is one grain of sand in the desert? analyzing individual neurons in deep NLP models. *AAAI*, 2019. doi: https://doi.org/10.1609/aaai.v33i01.33016309.

Guy Dar, Mor Geva, Ankit Gupta, and Jonathan Berant. Analyzing transformers in embedding space. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 16124–16170, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.186 53/v1/2023.acl-long.893. URL https://aclanthology.org/2023.acl-long.893.

Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. LLM.int8(): 8-bit matrix multiplication for Transformers at scale. *NeurIPS*, 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/c3ba4962c05c49636d4c6206a97e9c8a-Paper-Conference.pdf.

Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathematical framework for transformer circuits, 2021. URL https://transformer-circuits.pub/2021/framework/index.html.

Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition, 2022. URL https://transformer-circuits.pub/2022/toy_model/index.html.

Amit Elhelo and Mor Geva. Inferring functionality of attention heads from their parameters, 2024. URL https://arxiv.org/abs/2412.11965.

Kawin Ethayarajh. How contextual are contextualized word representations? Comparing the geometry of BERT, ELMo, and GPT-2 embeddings. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 55–65, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1006. URL https://aclanthology.org/D19-1006.

Team Gemma. Gemma. 2024. doi: 10.34740/KAGGLE/M/3301. URL https://www.kaggle.com/m/3301.

Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value memories. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wentau 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.

Mor Geva, Avi Caciularu, Kevin Wang, and Yoav Goldberg. Transformer feed-forward layers build predictions by promoting concepts in the vocabulary space. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 30–45, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.3. URL https://aclanthology.org/2022.emnlp-main.3.

Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. Dissecting recall of factual associations in auto-regressive language models. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 12216–12235, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.751. URL https://aclanthology.org/2023.emnlp-main.751.

Dirk Groeneveld, Iz Beltagy, Evan Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, Ananya Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, Shane Arora, David Atkinson, Russell Authur, Khyathi Chandu, Arman Cohan, Jennifer Dumas, Yanai Elazar, Yuling Gu, Jack Hessel, Tushar Khot, William Merrill, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew Peters, Valentina Pyatkin, Abhilasha Ravichander, Dustin Schwenk, Saurabh Shah, William Smith, Emma Strubell, Nishant Subramani, Mitchell Wortsman, Pradeep Dasigi, Nathan Lambert, Kyle Richardson, Luke Zettlemoyer, Jesse Dodge, Kyle Lo, Luca Soldaini, Noah Smith, and Hannaneh Hajishirzi. OLMo: Accelerating the science of language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15789–15809, Bangkok,

- Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-lon g.841. URL https://aclanthology.org/2024.acl-long.841.
 - Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. Finding neurons in a haystack: Case studies with sparse probing, 2023. URL https://arxiv.org/pdf/2305.01610.
 - Wes Gurnee, Theo Horsley, Zifan Carl Guo, Tara Rezaei Kheirkhah, Qinyi Sun, Will Hathaway, Neel Nanda, and Dimitris Bertsimas. Universal neurons in gpt2 language models, 2024. URL https://arxiv.org/pdf/2401.12181.
 - Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL https://arxiv.org/abs/2310.06825.
 - Olga Kovaleva, Saurabh Kulshreshtha, Anna Rogers, and Anna Rumshisky. BERT busters: Outlier dimensions that disrupt transformers. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 3392–3405, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.fin dings-acl.300. URL https://aclanthology.org/2021.findings-acl.300.
 - Vedang Lad, Wes Gurnee, and Max Tegmark. The remarkable robustness of llms: Stages of inference?, 2024. URL https://arxiv.org/abs/2406.19384.
 - 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.
 - Beren Millidge and Sid Black. The singular value decompositions of transformer weight matrices are highly interpretable, 2022. URL https://www.lesswrong.com/posts/mkbGjzxD8d8XqKHzA/the-singular-value-decompositions-of-transformer-weight.
 - Ari S. Morcos, David G.T. Barrett, Neil C. Rabinowitz, and Matthew Botvinick. On the importance of single directions for generalization, 2018. URL https://arxiv.org/pdf/1803.06959.pdf.
 - Neel Nanda. Neuroscope. Website, 2022. URL https://neuroscope.io.
 - Neel Nanda and Joseph Bloom. Transformerlens. https://github.com/TransformerLensOrg/TransformerLens, 2022.
 - Jingcheng Niu, Andrew Liu, Zining Zu, and Gerald Penn. What does the knowledge neuron thesis have to do with knowledge?, 2024. URL https://arxiv.org/pdf/2405.02421.
 - nostalgebraist. Interpreting gpt: The logit lens, 2020. URL https://www.lesswrong.com/ posts/Ackrb8wDpdaN6v6ru/interpreting-gpt-the-logit-lens.
 - Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry of large language models. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 39643–39666. PMLR, 21–27 Jul 2024. URL https://proceedings.mlr.press/v235/park24c.html.
 - Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.
 - Cody Rushing and Neel Nanda. Explorations of self-repair in language models. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 42836–42855. PMLR, 21–27 Jul 2024. URL https://proceedings.mlr.press/v235/rushing24a.html.

Noam Shazeer. Glu variants improve transformer, 2020. URL https://arxiv.org/pdf/2002.05202.

Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, Valentin Hofmann, Ananya Jha, Sachin Kumar, Li Lucy, Xinxi Lyu, Nathan Lambert, Ian Magnusson, Jacob Morrison, Niklas Muennighoff, Aakanksha Naik, Crystal Nam, Matthew Peters, Abhilasha Ravichander, Kyle Richardson, Zejiang Shen, Emma Strubell, Nishant Subramani, Oyvind Tafjord, Evan Walsh, Luke Zettlemoyer, Noah Smith, Hannaneh Hajishirzi, Iz Beltagy, Dirk Groeneveld, Jesse Dodge, and Kyle Lo. Dolma: an open corpus of three trillion tokens for language model pretraining research. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15725–15788, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.840. URL https://aclanthology.org/2024.acl-long.840.

Alessandro Stolfo, Ben Wu, Wes Gurnee, Yonatan Belinkov, Xingyi Song, Mrinmaya Sachan, and Neel Nanda. Confidence regulation neurons in language models, 2024. URL https://arxiv.org/abs/2406.16254.

Mingjie Sun, Xinlei Chen, J. Zico Kolter, and Zhuang Liu. Massive activations in large language models. *arXiv*, 2024. URL https://arxiv.org/pdf/2402.17762.

William Timkey and Marten van Schijndel. All bark and no bite: Rogue dimensions in transformer language models obscure representational quality. 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. 4527–4546, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.e mnlp-main.372. URL https://aclanthology.org/2021.emnlp-main.372.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *ArXiv*, abs/2302.13971, 2023. URL https://arxiv.org/pdf/2302.13971.

Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Neural Information Processing Systems*, 2017. URL https://api.semanticscholar.org/CorpusID:13756489.

Roman Vershynin. *High dimensional probability*. 2nd edition, 2025. URL https://www.math.uci.edu/~rvershyn/papers/HDP-book/HDP-2.pdf.

Elena Voita, Javier Ferrando, and Christoforos Nalmpantis. Neurons in large language models: Dead, n-gram, positional. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 1288–1301, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.75. URL https://aclanthology.org/2024.findings-acl.75.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2 technical report. arXiv preprint arXiv:2407.10671, 2024.

A LIMITATIONS

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

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, appendix E).

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.

The most expensive part was the activation-based analysis in appendix F.1: We needed a single run of ≈ 25 h to store the max/min activating examples for all neurons, and then ≈ 45 s per neuron (≈ 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 (section 2.3 and appendix C). 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.

C RANDOM BASELINES

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

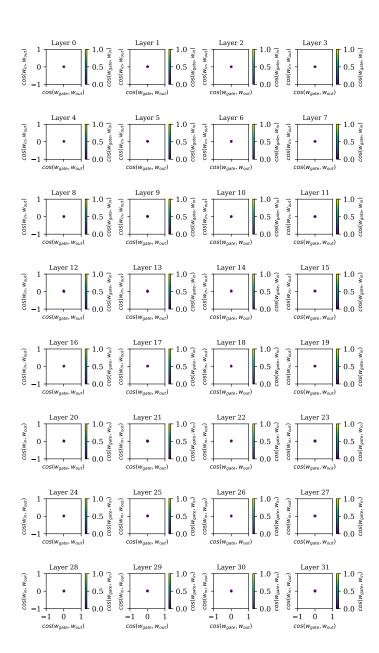


Figure 7: Equivalent of figure 3 for the randomly initialized OLMo-7B model (training checkpoint 0). Whatever doesn't look like this, is significant.

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

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

D ABLATION EXPERIMENTS

D.1 HYPOTHESES AND CHOICE OF METRICS

We originally had two hypotheses (which turned out to be wrong, see section 4):

- We hypothesized that *conditional strengthening* neurons might contribute to *subject enrichment* Geva et al. (2023), a crucial step of factual recall that involves MLPs writing appropriate attributes for the given subject. Both phenomena occur in roughly the same layers, and similar w_{in} and w_{out} could correspond to related concepts.
- We expected that *weakening* neurons would make the output distribution flatter, i.e. *increase* the entropy. This could happen by reducing the probability of high-ranking tokens (weakening directions corresponding to tokens) or by increasing the probability of very low-ranking tokens (weakening directions corresponding to negations of tokens).

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

D.2 Number of neurons to ablate

In preliminary experiments, we tried ablating 24 or 243 neurons in each run, which is the number of strengthening or weakening neurons in OLMo, respectively. Ablating 24 neurons (of any class) did not have a clear impact on any metric (including loss), so we sticked 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.

D.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.⁷ To improve replicability, we publish our complete code as well as our subject-attribute dataset.

⁶https://stats.stackexchange.com/questions/85916/distribution-of-scalar-products-of-two-random-unit-vectors-in-d-dimensions

⁷A 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 appendix F.1

Neuron	RW category	$\cos(oldsymbol{w}_{ ext{gate}}, oldsymbol{w}_{ ext{in}})$	$\cos(oldsymbol{w}_{ ext{gate}}, oldsymbol{w}_{ ext{out}})$	$\cos(oldsymbol{w}_{ m in}, oldsymbol{w}_{ m out})$
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

E DOUBLE CHECKING

In our case studies, we observe that many neurons have the property of **double checking**: The two reading weight vectors (w_{gate} and w_{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 w_{gate} and w_{in} have a high overlap. More formally, let $U = \{u_0, ..., u_{d_{\text{vocab}}}\}$ be the set of unembedding vectors; then

$$\underset{u \in U}{\arg \max} \cos(u, \boldsymbol{w}_{\text{gate}}) \approx \underset{u \in U}{\arg \max} \cos(u, \boldsymbol{w}_{\text{in}}).$$

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 $\boldsymbol{w}_{\text{gate}}, \boldsymbol{w}_{\text{in}}$ that are reasonably similar to a u_i but not to any other u_j . These $\boldsymbol{w}_{\text{gate}}, \boldsymbol{w}_{\text{in}}$ can even be (approximately) orthogonal to each other, as in the following three-dimensional toy example: $u_1 = (1,0,0), u_2 = (0,1,0), \boldsymbol{w}_{\text{gate}} = (1,0,1), \boldsymbol{w}_{\text{in}} = (1,0,-1).$

However the phenomenon is *unlikely* to occur in random vectors, and hence is a significant finding: If choosing w_{gate} , w_{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 w_{gate} and w_{in} is to detect a concept (e.g. a token prediction) represented by a vector u, then we would get the best performance with $w_{\text{gate}} = w_{\text{in}} = u$, i.e., w_{gate} , w_{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.

F CASE STUDIES

F.1 CASE STUDY OF A WEAKENING NEURON

We now qualitatively examine a weakening neuron in more detail, showing that these neurons can have a quite complex behavior. In appendix F we detail the methods, including how we chose the neuron, and present many more case studies from various RW functionalities.

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

 $^{^8}$ We separately publish code and data for the second method at https://anonymous.4open.science/r/neuroscope-B446

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)

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

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).

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

We choose neuron **31.9634**. It is also a *prediction* neuron in the sense of (Gurnee et al., 2024), which indicates that it directly promotes a specific set of tokens.

Our analysis can be summarized as follows: This neuron produces "again" when the residual stream contains "minus again"; but the examples activating the neuron do not have an obvious semantic relationship to again.

More specifically, the weight-based analysis indicates that w_{out} of this neuron is closest to forms of again. Judging by the weights, the neuron activates negatively when the residual stream contains information both for and against predicting again, and then weakens the again direction. It activates positively when the residual stream contains the "minus again" direction, and then weakens that direction.

Surprisingly, despite its strong positive cosine similarity $(\cos(\mathbf{w}_{\text{gate}}, \mathbf{w}_{\text{in}}) = 0.7164)$, the neuron often activates negatively. On the negative side, strong activations are often on punctuation, and the actual next token is often *meanwhile* or *instead* (instead of *again*). The neuron may ensure only these tokens are predicted, and not the relatively similar *again*. On the positive side, the activations do not have any obvious semantic relationship to *again*. 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. There are also weaker positive activations when *again* is a plausible continuation, e.g., on the token *once*. In these cases, *again* is already weakly present in the residual stream before the last MLP. Accordingly, $Swish(\mathbf{w}_{\text{gate}} \cdot \mathbf{x}_{\text{norm}})$ is weakly negative (but distinct from 0), and $\mathbf{w}_{\text{in}} \cdot \mathbf{x}_{\text{norm}} < 0$, which leads to a positive activation and thus reinforces *again*.

 $^{^{9}}$ The notation is "layer.neuron", with zero-based indexing. The model – still OLMo-7B – has 32 layers, so this neuron is in the final layer.

In contrast to this neuron, **strengthening neuron 28.4737** is also a prediction neuron (it also directly promotes a specific set of tokens), but its read-write behavior is much more straightforward: It further promotes *review* when the residual stream already indicates that this is the obvious next token.

These two case studies show that even when the output weights are highly interpretable, strengthening and weakening have a very different overall behavior, and the weakening behavior is much more complex.

F.2 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 $\cos(W_U, w_{\text{out}})$ 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 $\cos(\boldsymbol{w}_{in}, \boldsymbol{w}_{out})$ among those neurons whose $\cos(\boldsymbol{w}_{gate}, \boldsymbol{w}_{out})$ is within the randomness range (section 2.3 and appendix C). 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).

F.3 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.

F.4 RESULTS AND ANALYSIS FOR PREDICTION NEURONS

See table 3.

Strengthening neuron 28.4737¹⁰ 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_{gate} and w_{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 appendix F.1.

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.)

¹⁰The notation is "layer.neuron", with zero-based indexing.

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 $(-w_{in})$. An example for a large negative activation is *puts you on multiple webpages* at. Conversely, if *when* is absent, and time-related tokens are absent too, the neuron activates positively and suppresses *when* further.

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

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

F.5 RESULTS AND ANALYSIS FOR PROTOTYPICAL EXAMPLES

See table 4.

Strengthening neuron 25.9997.

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

Weakening neuron 31.7117.

Conditional weakening neurons 23.6543.

Proportional change neuron 25.7415.

F.6 MORE CASE STUDIES

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

F.6.1 CONDITIONAL STRENGTHENING NEURONS

- **0.1480**: $w_{\text{gate}}, -w_{\text{in}}, -w_{\text{out}}$ all have tokens similar to box (when using W_E). Activates on Xbox.
- **4.1940**: country appears in w_{in} among many other things. When using W_E , Philippines and Manila appear in w_{out} . Activates on Philippines.
- **4.3720**: gate seems country/government related. When using W_E , we find w_{out} , w_{gate} contain some country names. Activates on *Denmark*.
- 4.4801: Muhammad appears in the gate vector. Activates on Muhammad.

¹¹The actual sentence ends with *as soon as* and comes from a now-dead webpage. We also found one occurrence of *at when* in what seems to be a paraphrase of the same text, on https://www.docdroid.net/RgxdG5s/fantastic-tips-for-bloggers-of-all-amountsoystcpdf-pdf. We suspect that 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*.

- 4.5772: predicts ian as in Egyptian. When using W_E , all three weight vectors contain Egypt. Activates on Egypt.
- 4.6517 has a very Ireland (or Celtic nations) related gate vector. The interpretations of the other two weights are less obvious, but *Irish* and *Dublin* appear in w_{in} among many other things, and UK and *London* appear in $-w_{out}$ (Ireland is emphatically *not* in the UK!) When using W_E , *Ireland* appears among the top tokens of all three weight vectors. Activates on *Ireland*.
- 4.6799: When using W_E , Vietnam is among the tokens corresponding to $-w_{\text{out}}$. Activates on Vietnam
- **4.7667**: all three weights related to consoles in different ways. Activates on *Xbox*
- 4.9983: w_{out} is related to electronic devices, w_{in} either electronic devices or sports (surfing may belong to both), w_{gate} is also mostly related to electronic devices. When using W_E , we find w_{out} contains iPhone as a top token. Activates on iPhone.
- **4.10859**: When using W_E , we find w_{gate} , w_{out} include *Thailand* as a top token, w_{out} additionally *Buddha*, *Buddhist*. Activates on *Thailand*.
 - **4.10882**: When using W_E , we find $-\mathbf{w}_{out}$ contains Italy, $-\mathbf{w}_{in}$, \mathbf{w}_{gate} additionally contain Rome. Activates on Italy.
 - **4.10995**: Boston appears in gate and Massachusetts in $-\mathbf{w}_{in}$. When using W_E , we find $-\mathbf{w}_{out}$, \mathbf{w}_{gate} contain Massachusetts and Boston, $-\mathbf{w}_{in}$ contains Boston. Activates on Massachusetts.
- 22.2589: w_{gate} and $-w_{\text{in}}$ recognize tokens like *Islam*, *Muhammad* and others related to the Arabo-Islamic world. The same goes for $-w_{\text{out}}$ (as it is similar to w_{in}). Activates on *Muhammad*.
 - **24.4880**: For all three weight vectors the first four tokens (but not more) are Philippine-related (even though the gate vector is actually not very similar to the others). The gate vector also reacts to other geographical names, which *may* have in common that they are associated with non-"white" (Black, Asian or Latin) people in the US sense (*Singapore, Malaysian, Nigerian, Seoul, Pacific, Kerala, Bangkok*, but also (*Los) Angeles* and *Bronx*). Activates on *Philippines*.
 - **24.6771**: w_{gate} , $-w_{\text{in}}$, $-w_{\text{out}}$ all correspond to capitalized first names. Activates on *Muhammad*.
- 25.2723: Some tokens associated with w_{in} and w_{out} are possible completions for th (th-ousand, th-ought, th-orn. When using W_E , in all three weights there are a few th tokens, but also with ph and similar. Activates on Thailand.
 - **25.10496**: $-w_{\text{in}}$, $-w_{\text{out}}$ correspond to tokens starting with v (upper or lower case, with or without preceding space). w_{gate} on the other hand seems to react to appropriate endings for tokens starting in v: vol-atility, v-antage, v-intage, v-locity, V-ancouver. When using W_E , we also find all three weight-vectors are very v-heavy. Activates on V-ietnam.

F.6.2 Weakening neurons

1095

1099

1102

1103

1104

1105

1106 1107

1108

1112

1113

1114

111511161117

1118

1119

112011211122

11231124

1125

1126

30.9996: Downgrades weird tokens if present / promotes frequent English stopwords if absent. Also an attention deactivation neuron for 15 heads in layer 31.

F.6.3 Proportional Change Neurons

- **25.7032**: Some tokens associated with w_{gate} and w_{out} are possible completions for x or ex (X-avier, x-yz, ex-cel, ex-ercise. When using W_E , both x and box (with variants) appear in all three weight vectors. Activates on Xbox.
- 25.8607: All three vectors correspond to tokens related to cities. Moreover, $-w_{out}$ seems to correspond to non-city places, such as national governments or villages. w_{in} is actually not that similar to w_{gate} , w_{out} (in terms of cosine similarities), but all three correspond to city-related tokens. When using W_E , in all three weights there are a few city-related tokens. Activates on *Paris*. We may think of the two input directions as two largely independent ways of checking that "it's about a city" (this is a recurring phenomenon that we describe in appendix E). When the gate activates but the linear input does not confirm it's about a city, the output promotes closely related but non-city interpretations (for example *Paris* actually refers to the French government in some contexts).

- 29.8118: Partition neuron, highest variance of all proportional change neurons. Also an attention deactivation neuron for 4 heads (0,2,11,15) in layer 30.
- 31.5490: Activates on *Muhammad*. w_{gate} reacts to various Asian names and Asian-sounding subwords, w_{in} to surnames as opposed to other English words starting with space and uppercase letter. w_{out} corresponds to more Asian stuff (mostly subwords) as opposed to English surnames.
 - 31.6275: Mostly promotes two-letter tokens (no preceding space, typically uppercase). $-w_{in}$ typically lowercase single letters. $-w_{gate}$ mostly lowercase two-letter tokens. "If no lowercase two-letter tokens, promote uppercase two-letter tokens proportionally to absence of lowercase single letters"?
- 31.8342: This is an -ot- neuron: w_{gate} and w_{out} correspond to -o(t)- suffixes, $-w_{\text{in}}$ to various -ot- stuff.

 Judging by the weight similarities, we expect that w_{out} is typically activated negatively: downgrade -o(t)- suffixes if present in the residual stream. Activates on Egypt.

1140

1141

1142

1146

1148 1149

1150

1151

1161

1162

1163

1164

1165

1166

1167

1168

1169

1172

1173 1174

1175

1176

F.6.4 ORTHOGONAL OUTPUT NEURONS

- **0.1758**: When using W_E , all three weight vectors' top tokens are famous web sites, including *YouTube*. Activates on *YouTube*.
- **0.3338**: When using W_E , we find especially w_{gate} and $-w_{\text{in}}$, but also $-w_{\text{out}}$ are similar to smartphone-related tokens. Activates on *iPhone*.
- 0.3872: When using W_E , we find especially w_{gate} , but also $-w_{\text{in}}$ and $-w_{\text{out}}$ correspond to city names. Activates on *Paris*.
- 0.7829: When using W_E , we find w_{in} , w_{out} and to a lesser extent w_{gate} correspond in large part to software names. Activates on *iTunes*.
- **0.7966**: When using W_E , the weight vectors mostly correspond to tokens starting with *th*. Activates on *Thor*.
 - **29.2568**: $\boldsymbol{w}_{\text{out}}$ Asian (Thai?) sounding syllables vs. (Asian) geographic names in English and other stuff; $\boldsymbol{w}_{\text{in}}$ reacts to Thailand and Asian (geography) stuff as opposed to (mostly) US stuff; $\boldsymbol{w}_{\text{gate}}$ pretty much the same. Activates on *Thailand*.
 - **29.3327**: w_{gate} mostly reacts to city names (*Paris* being the most important one), $-w_{\text{in}}$ countries and cities, especially in continental Europe (*France* and *Paris* on top) as opposed to stuff related to the former British Empire. Relevant is $-w_{\text{out}}$ which corresponds to pieces of geographical names and especially rivers in France (*Se-ine, Rh-one / Rh-ine, Mar-ne, Mos-elle... Norm-andie, Nancy, commun...*). w_{gate} and $-w_{\text{in}}$ also react to viver(s). Activates on viver(s).
- 29.4101: w_{gate} and w_{in} react to *YouTube* (top token!), w_{out} downgrades it (almost bottom token) and promotes *subscrib**, *views*, *channels* etc. Activates on *YouTube*.
 - **29.6417**: Downgrades *recording* and similar. w_{gate} and w_{in} are also similar and involve *iTunes*. Activates on *iTunes*.
 - **29.9734**: w_{gate} reacts to the East in a broad sense as opposed to the West (*Iran, Kaz-akhstan, Kash-mir, Ukraine...*), w_{in} mostly to male first names without preceding space. w_{out} seems to produce word pieces that could begin a foreign name. Activates on *Muhammad*.
- 30.2667: w_{gate} reacts to suffixes (for adjectives derived from place names) like *en*, *ian*, *ians*, basically the same for w_{in} and w_{out} . Activates on *Muhammad*.
- 30.3143: w_{gate} reacts to words related to entities that are authoritative for various reasons (officials, authorities, according, researchers, spokesman, investigators...). $-w_{\text{in}}$ reacts to uncertainty (reportedly, according... allegedly... accused). $-w_{\text{out}}$ is again police, authorities, officials, court but with no preceding space. Activates on Philippines. What authorities and uncertainty have to do with the Philippines is unclear.
- 30.3883: w_{gate} and $-w_{\text{in}}$ react to *Virginia* and *Afghanistan*, among others (in the case of w_{gate} : as opposed to other geographical names with no preceding space associated with the South and the sea); $-w_{\text{out}}$ is activated and promotes all variants of af (and ghan) but downgrades Virginia etc. Activates on *Afghanistan*.

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

30.4577: Seems to be related to rugby: w_{gate} and slightly less obviously w_{in} react to rugby-related tokens (*midfielder*, *quarterback*...); w_{out} promotes different tokens that upon reflection could be related to rugby as well. Activates on *Ireland*.

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

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

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

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

G MORE PLOTS

These final figures show additional results:

- Table 5 and figures 8 to 10 show more results on activation frequencies: by discrete classes, on all layers, and plotted against $\cos(w_{\text{gate}}, w_{\text{in}})$ or $\cos(w_{\text{gate}}, w_{\text{out}})$.
- Figures 11 to 13 show that ablating 243 neurons from other classes than weakening does not have any effect on attribute rate.
- Figures 14 to 21 show the effect of various ablations on entropy, loss, rank and scale.
- From figure 22, we show our analyses of RW functionalities by layer (section 3) for all the models we investigated.

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

- In Yi and the OLMo models, the prevalence of conditional strengthening neurons starts even earlier, at the very first layer. A particularly interesting example is Yi: In layer 0 an enormous 68% of all neurons are conditional strengthening, then almost none, then there is a second wave around layers 11-17 (out of 32) which have around 25% of conditional strengthening neurons each.
- In some models, especially the OLMo ones, there is a non-negligible number of conditional weakening neurons. They tend to appear in middle-to-late layers, shortly after the conditional

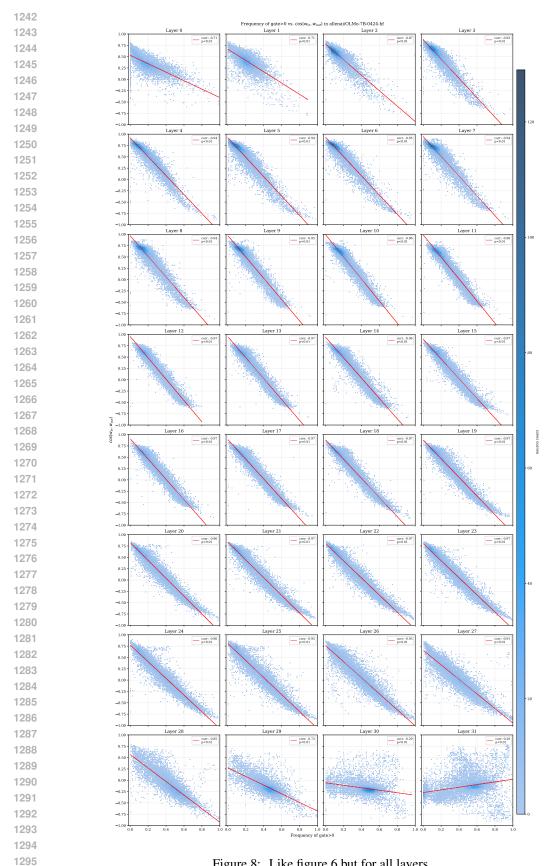


Figure 8: Like figure 6 but for all layers.

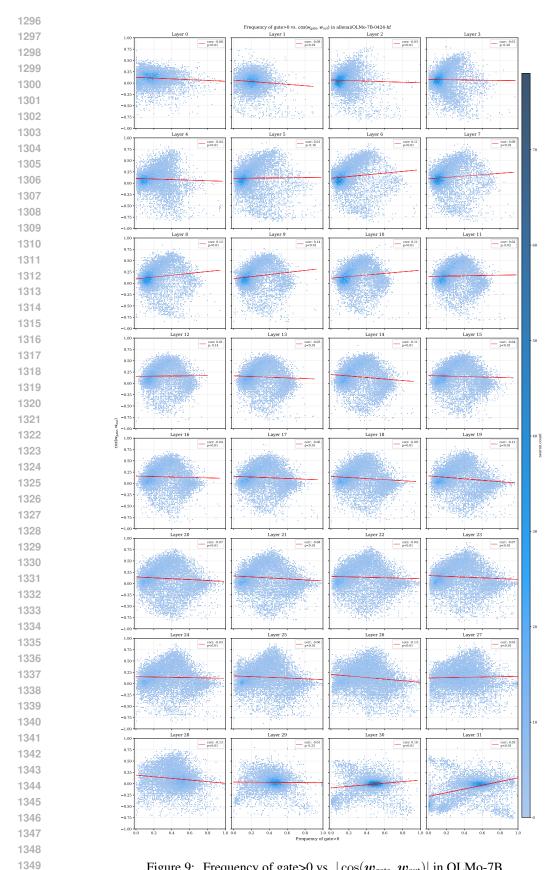


Figure 9: Frequency of gate>0 vs. $|\cos(w_{\text{gate}}, w_{\text{out}})|$ in OLMo-7B.

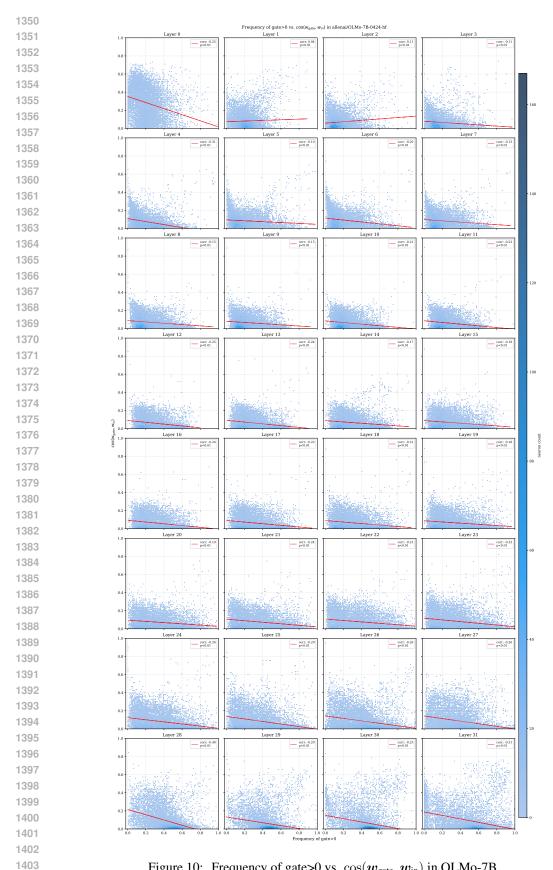


Figure 10: Frequency of gate>0 vs. $\cos(\boldsymbol{w}_{\text{gate}}, \boldsymbol{w}_{\text{in}})$ in OLMo-7B.

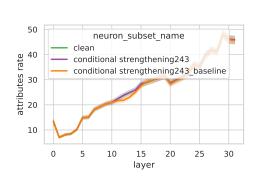


Figure 11: Attribute rate when ablating 243 conditional strengthening neurons.

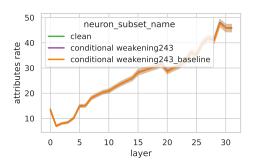


Figure 12: Attribute rate when ablating 243 conditional weakening neurons.

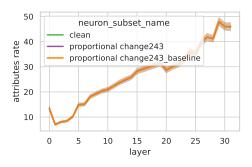


Figure 13: Attribute rate when ablating 243 proportional change neurons.

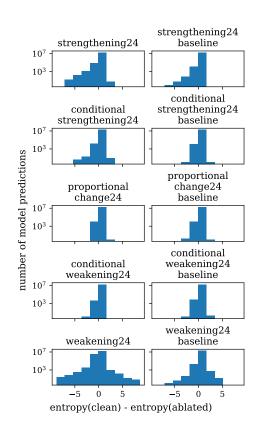


Figure 14: Effect on entropy when ablating 24 neurons from various RW classes.

strengthening wave. The clearest example is OLMo-1B, with a peak of 1418 conditional weakening neurons out of 8192 (17%) in layer 9 out of 16.

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

- For almost all neurons the cosine similarities are still clearly below 1 (the dots do not fill out the edges in figure 3). This echoes and extends Gurnee et al. (2024)'s findings that in GPT2 the IO cosine similarity is approximately bounded by ±0.8. In other words, we almost never get the *prototypical* cases of conditional strengthening / weakening etc., as defined in section 2. This might be an effect of randomness (strong cosine similarities are less likely), but could also suggest that even input manipulator neurons add some novel information to the residual stream.
- We also observe that for the vast majority of neurons, cos(w_{gate}, w_{in}) ≈ 0: This can be seen
 in the boxplots in the appendix, as well as the purple color in figure 3. Thus most neurons
 operate on two input directions in the residual stream (not a single one), resulting in higher
 expressivity and more complex semantics. If not random, this could be related to double
 checking; see appendix E.

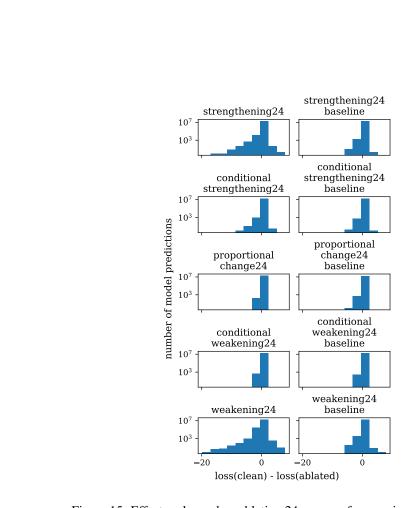


Figure 15: Effect on loss when ablating 24 neurons from various RW classes.

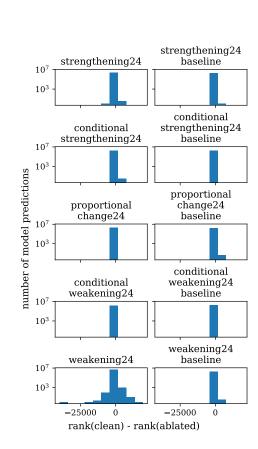


Figure 16: Effect on correct token rank when ablating 24 neurons from various RW classes.

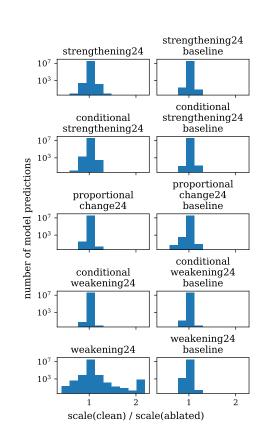


Figure 17: Effect on scale of last hidden state when ablating 24 neurons from various RW classes.

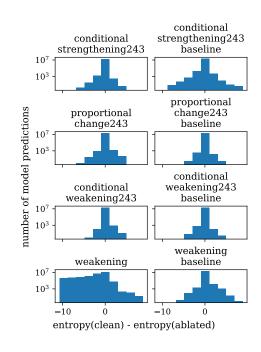


Figure 18: Effect on entropy when ablating 243 neurons from various RW classes.

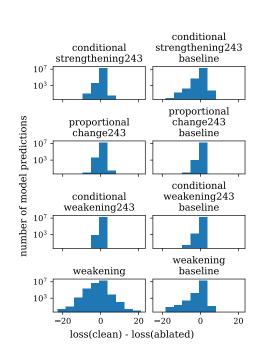


Figure 19: Effect on loss when ablating 243 neurons from various RW classes.

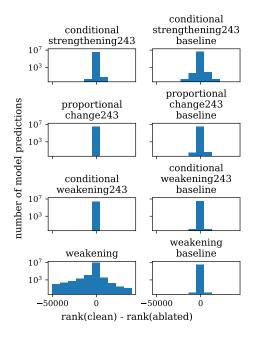


Figure 20: Effect on correct token rank when ablating 243 neurons from various RW classes.

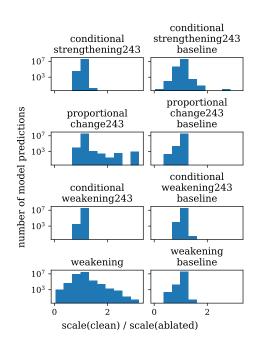


Figure 21: Effect on scale of last hidden state when ablating 243 neurons from various RW classes.

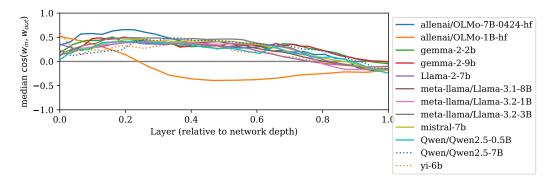


Figure 22: Median of $\cos(w_{\rm in}, w_{\rm out})$ by layer (x-axis) for all 12 models investigated. Unlike figure 1 we also include the models of 1B parameters and below. All models follow the same general pattern, but OLMo-1B switches to negative values earlier than the others.

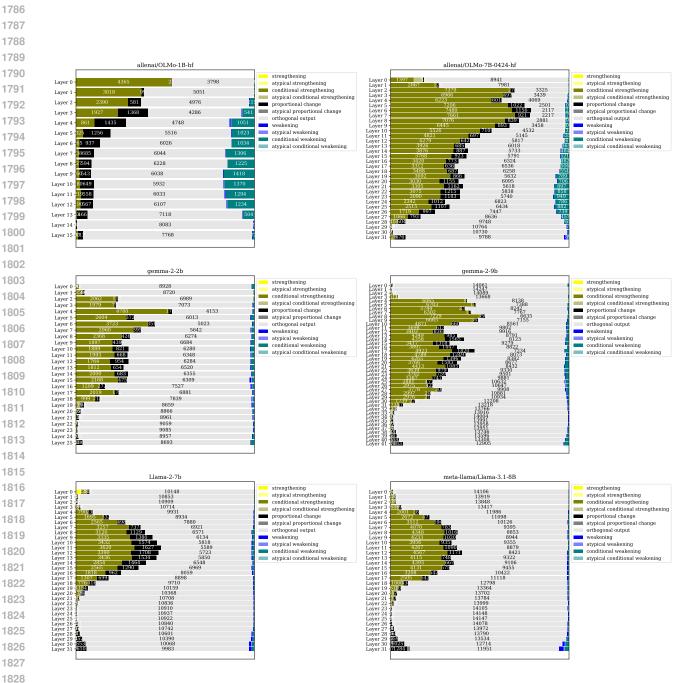


Figure 23: Distribution of neurons by layer and category for a range of models



Figure 24: Continuation of figure 23. Including a copy of figure 2 (Llama-3.2-3B) for convenience.

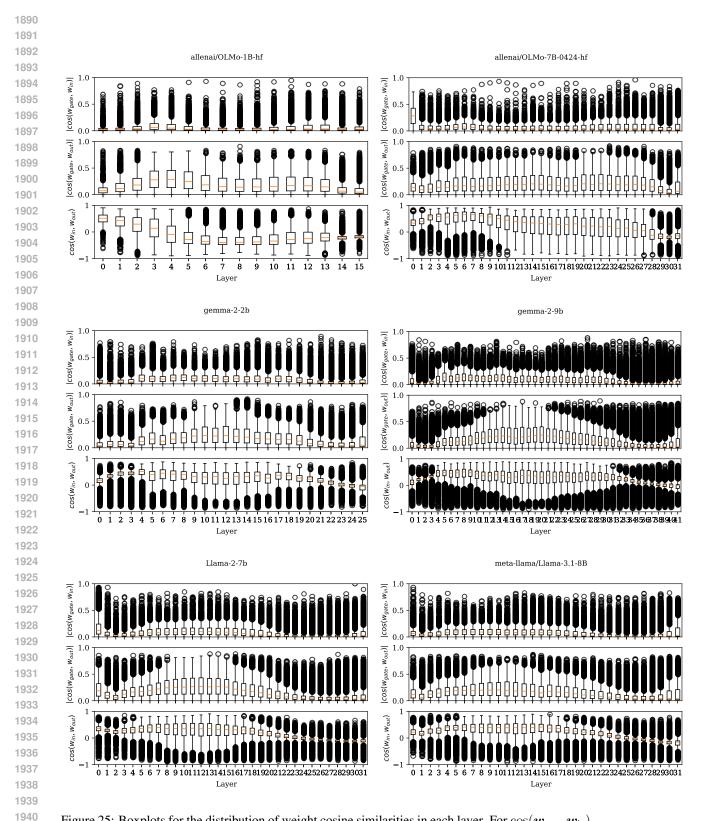


Figure 25: Boxplots for the distribution of weight cosine similarities in each layer. For $\cos(w_{\text{gate}}, w_{\text{in}})$ and $\cos(w_{\text{gate}}, w_{\text{out}})$ we show the absolute value since their sign does not carry any information on its own.

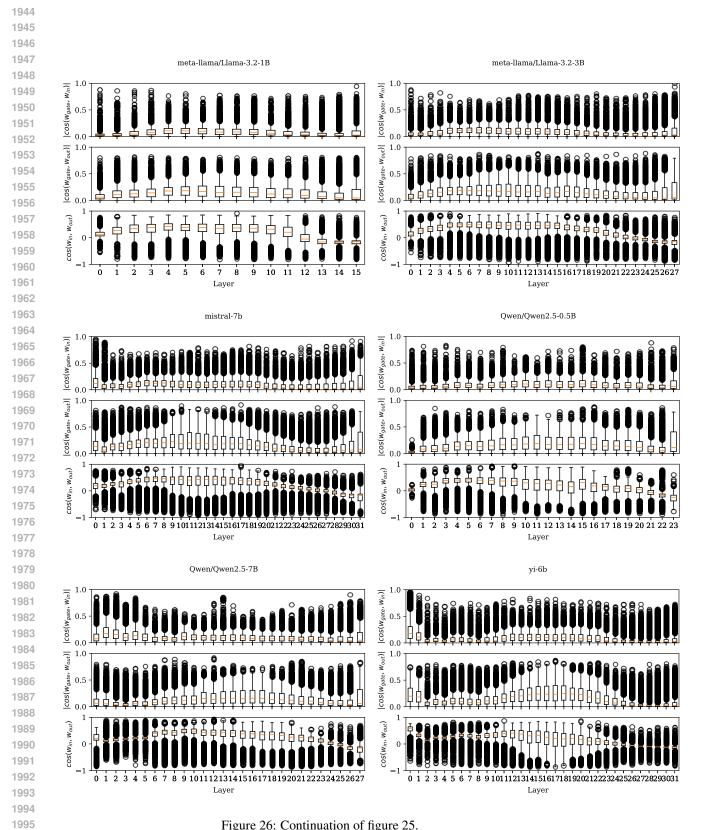


Figure 26: Continuation of figure 25.

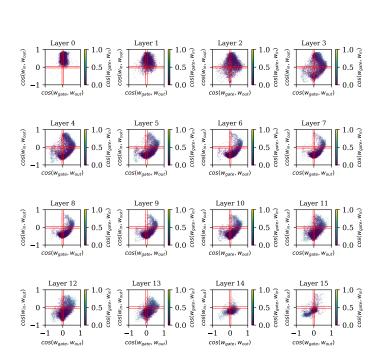


Figure 27: Equivalent of figure 3 for OLMo-1B

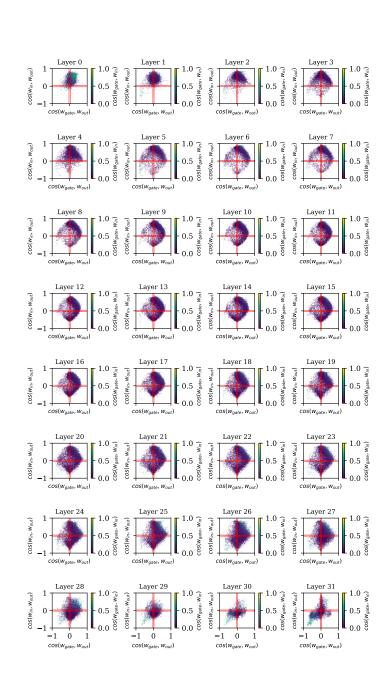


Figure 28: Equivalent of figure 3 for OLMo-7B-0424

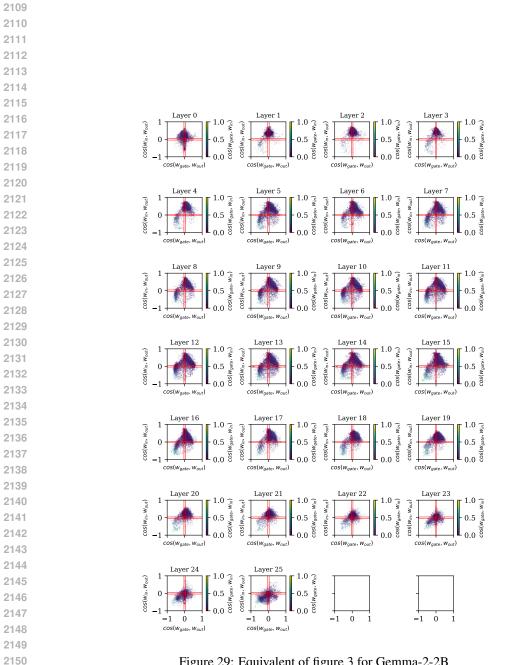


Figure 29: Equivalent of figure 3 for Gemma-2-2B

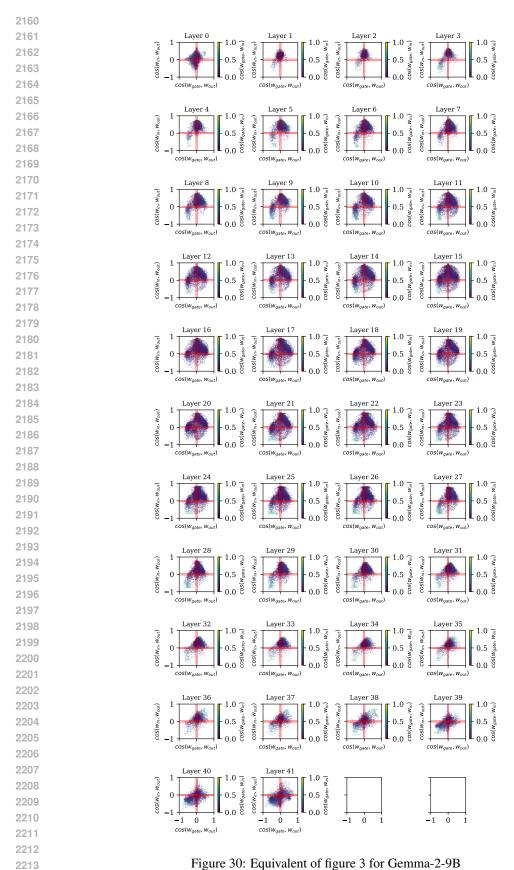


Figure 30: Equivalent of figure 3 for Gemma-2-9B

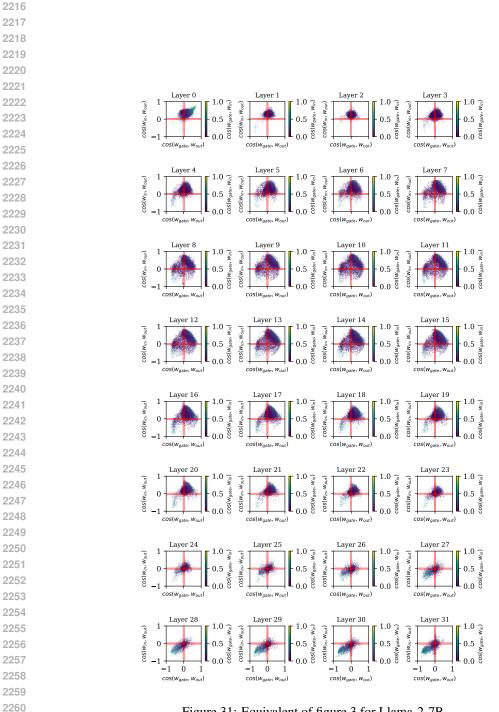


Figure 31: Equivalent of figure 3 for Llama-2-7B

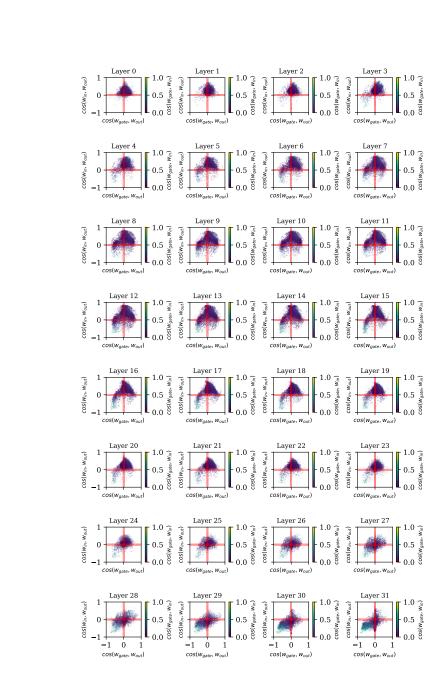


Figure 32: Equivalent of figure 3 for Llama-3.1-8B

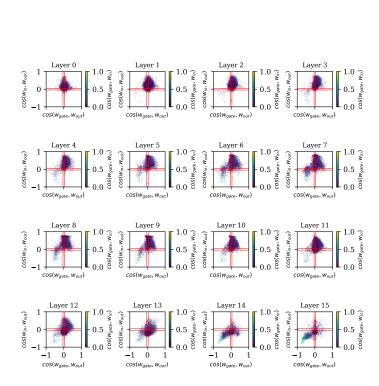


Figure 33: Equivalent of figure 3 for Llama-3.2-1B

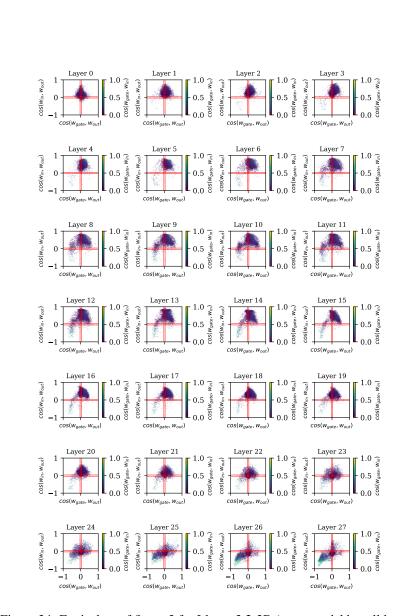


Figure 34: Equivalent of figure 3 for Llama-3.2-3B (same model but all layers)

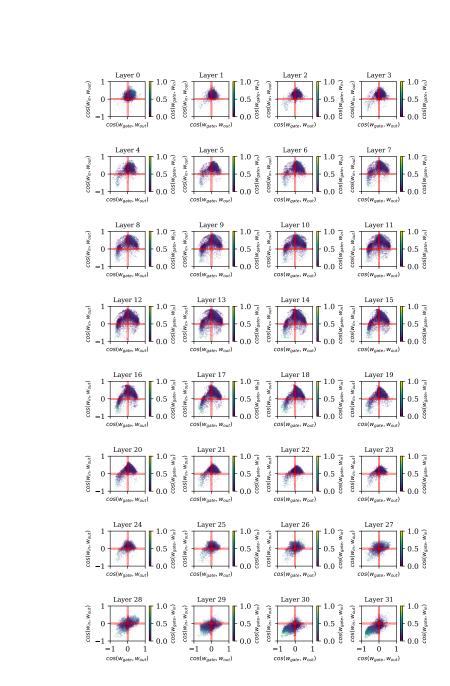


Figure 35: Equivalent of figure 3 for Mistral-7B

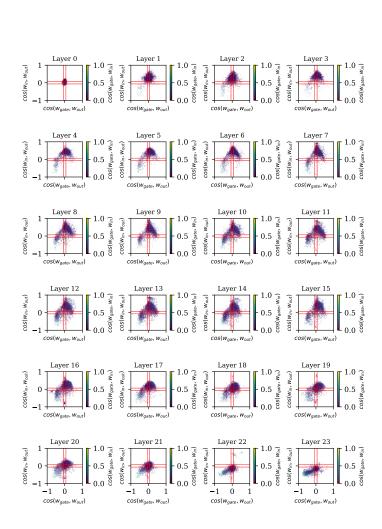


Figure 36: Equivalent of figure 3 for Qwen2.5-0.5B

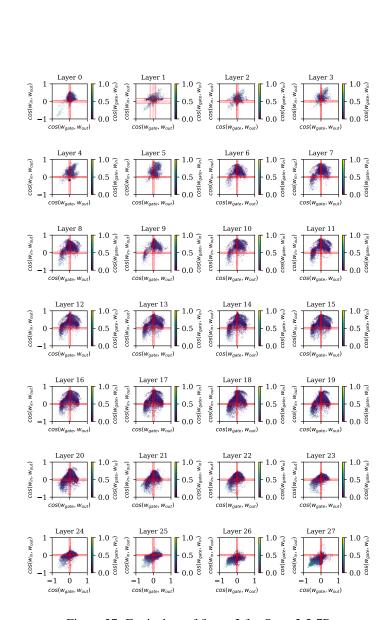


Figure 37: Equivalent of figure 3 for Qwen2.5-7B

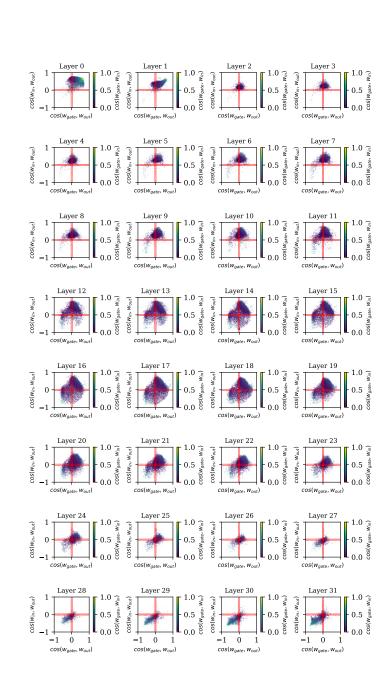


Figure 38: Equivalent of figure 3 for Yi-6B