# Auxiliary Metrics Help Decoding Skill Neurons in the Wild

**Yixiu Zhao**[*]  **Xiaozhi Wang**[*]  **Zijun Yao**  **Lei Hou**  **Juanzi Li**
Tsinghua University, Beijing, China, 100084
yx-zhao22@mails.tsinghua.edu.cn

## Abstract

Large language models (LLMs) exhibit remarkable capabilities across a wide range of tasks, yet their internal mechanisms remain largely opaque. In this paper, we introduce a simple, lightweight, and broadly applicable method with a focus on isolating neurons that encode specific skills. Building upon prior work that identified "skill neurons" via soft prompt training on classification tasks, our approach extends the analysis to complex scenarios involving multiple skills. We correlate neuron activations with auxiliary metrics—such as external labels and the model's own confidence score, thereby uncovering interpretable and task-specific behaviors without the need for manual token aggregation. We empirically validate our method on tasks spanning open-ended text generation and natural language inference, demonstrating its ability to detect neurons that not only drive known skills but also reveal previously unidentified shortcuts in arithmetic reasoning on BigBench.

## 1 Introduction

Large language models (LLMs) have become increasingly powerful. As these models grow in complexity, interpretability research becomes essential—to understand their inner workings and also steer them safely (Bereska and Gavves, 2024).

Interpretability research typically unfolds in two stages (Räuker et al., 2023). In the *observation* stage, researchers attribute specific model behaviors to substructures within the network. For instance, prior studies have shown that certain neurons encode specific skills, like sentiment detection or fact retrieval (Radford et al., 2017; Gurnee et al., 2023; Wang et al., 2022; Bills et al., 2023; Choi et al., 2024). In the *intervention* stage, these insights are taken further by modifying the weights or activations of the identified substructures, thereby causally influencing the model's behavior.
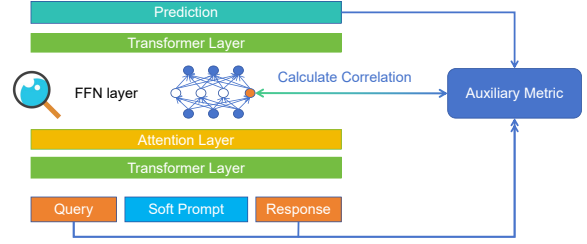
---

[*]indicates equal contribution.



Figure 1: **Overview of Our Methodology** We calculate the correlation between feedforward-layer neuron activations on a trained soft prompt and an auxiliary metric to identify skill-related neurons.

In this work, we introduce a simple, lightweight, and broadly applicable method for probing large language models during the observation phase, with a particular focus on identifying neurons associated with specific skills using a downstream corpus. Building on earlier work that isolated "skill neurons" through soft prompt training (Lester et al., 2021; Wang et al., 2022), our approach expands the analysis from single-skill classification tasks to more complex scenarios involving multiple skills. After training soft prompts on our target tasks, we leverage **auxiliary metrics**—such as auxiliary labels or the model's own confidence in its predictions to pinpoint neurons whose activations on the soft prompt are strongly correlated with these metrics corresponding to model skills.

Our framework is also closely related to network dissection, a technique widely used to uncover low-level features in computer vision. Seminal work by Bau et al. (2017) showed how individual neurons in convolutional networks can be linked to specific visual concepts by analyzing their activations in response to corresponding labels. In the NLP domain, similar methods have been employed to identify neurons that capture low-level linguistic phenomena. However, these approaches often struggle with complex, multi-token features because aggregating activations across tokens is

not straightforward (Gurnee et al., 2023). By calculating the activation on the trained soft prompt, we eliminate the need for manual aggregation, thereby better accommodating the inherent complexity of language. Moreover, soft prompt training taps into the full potential of the pretrained model, enabling us to uncover higher-level, task-specific skills encoded within the network.

We empirically demonstrate that our methods can effectively detect neurons with specific skill in open-ended generation (Skill-Mix, Yu et al. (2023)) and natural language inference (Heuristic Analysis for NLI Systems, McCoy et al. (2019)). Furthermore, the framework detects neurons corresponding to a previously unknown shortcut in the arithmetic subset in BigBench (Srivastava et al., 2023) without reliance on any auxiliary label.

## 2 Methodology

### 2.1 Preliminary

**Prompt Tuning.** Soft prompt tuning (Lester et al., 2021; Li and Liang, 2021) has been widely adopted to adapt language models to downstream tasks. For a pretrained decoder-based language model $\mathcal{M}$, given an input instruction $x$, the embedding function Emb of $\mathcal{M}$ will first maps these tokens to a sequence of $d$-dimensional vectors $\text{Emb}(x)$, the vectors will then be mapped to a distribution over the vocabulary set, which we will denote as $\text{Pr}_{\mathcal{M}}(\dots | \text{Emb}(x))$. Given a training set of instruction $x$ and desired completion $y$, prompt tuning will train $l$ randomly initialized vectors $\{p_1, \dots, p_l\}$ to minimize the following loss:

$$-\text{E}\left[\frac{1}{|y|}\sum_{j=1}^{|y|}\log \text{Pr}_{\mathcal{M}}(y_j \mid \text{Emb}(x), p_1, \dots p_l, \text{Emb}(y_{1:j-1}))\right].$$

We will interpret the pretrained model based on the neurons' activations on the trained prompt. In order to do so, we will (1) freeze the parameter of $\mathcal{M}$ and only train the prompt, and (2) put the soft prompt after the instruction such that the soft prompt can attend to the input instruction through the causal attention mechanism.

**Neurons.** Transformer (Vaswani et al., 2023) is the most common architecture used in language models. Each layer of a Transformer composed of one attention module and one feed-forward layers. Feed-forward layers in Transformers are typically a two-layer fully connected network operates independently on the hidden state of each token. In the standard LLaMA architecture, the feedforward layer of width $m$ is defined as,

$$\text{FFN}(h) = W^{(3)}\big((W^{(1)}h) \odot \text{SiLU}(W^{(2)}h)\big),$$

where $h \in \mathbb{R}^d$ is the input hidden state, $W^{(1)}, W^{(2)} \in \mathbb{R}^{m \times d}, W^{(3)} \in \mathbb{R}^{d \times m}$ are learnable parameters, $\odot$ denotes element-wise multiplication, and $\text{SiLU}(x) = \frac{x}{1+e^{-x}}$ is the SiLU activation function. We will follow the definition in Wang et al. (2022) to define the $i$-th neuron as the set of parameters $\{W_{i,:}^{(1)}, W_{i,:}^{(2)}, W_{:,i}^{(3)}\}$ for $i \in [m]$. The activation of the $i$-th neuron on hidden state $h$ is then defined as $W_{i,:}^{(1)}h\text{SiLU}(W_{i,:}^{(2)}h)$.

Because the input of the Transformer is always a sequence, each neuron will have different activations at different position of the sequence. We will focus on the activations on the trained soft prompt and use $a_{l,i,k}(x)$ to denote the activation on the $k$-th position of the soft prompt for neuron $i$ on layer $l$ when the input instruction is $x$.

### 2.2 Method

We introduces a systematic approach to identifying neurons in language models responsible for activating specific skills in response to tasks. Given a training set $S_{\text{train}}$ and a validation set $S_{\text{val}}$ consisting of instruction completion pairs. The method consists of the following stages:

1. **Training:** We will first train soft prompts using a frozen pretrained language model on $S_{\text{train}}$. This enables task-specific adaptation without altering the underlying model weights.

2. **Metric Calculation:** After training, we will select a *helper metric* $m$ for each sample in validation set, which is a function of the input sequence and the models' output that maps to a real number. For example, we can map all the sequence to their corresponding loss values given the trained model. We will use $m(S_{\text{val}})$ to denote the calculated metrics.

3. **Neuron Selection:** For every neuron $i$ on layer $l$, we will calculate their activations on the soft prompt over the instruction and completion of the validation set, we will denote the activations on the $k$-th position of the soft prompt as $a_{l,i,k}(S_{\text{val}})$. We then compute the Pearson correlation coefficient between $a_{l,i,k}(S_{\text{val}})$ and $m(S_{\text{val}})$, defined as:

$$\mathrm{corr}_{l,i,k} = \frac{\sum_{k=1}^{N}\left(a_{l,i,k} - \overline{a_{l,i,k}}\right)\left(m_k - \overline{m}\right)}{\sqrt{\sum_{k=1}^{N}\left(a_{l,i,k} - \overline{a_{l,i,k}}\right)^2}\sqrt{\sum_{k=1}^{N}\left(m_k - \overline{m}\right)^2}}$$

- $N$ is the number of samples in $S_{\mathrm{val}}$.
- $\overline{a_{l,i,k}}$ is the mean activation of neuron $i$ in layer $l$ on the $k$-th position of the soft prompt across the validation set.
- $\overline{m}$ is the mean of the helper metric across the validation set.

We will then define the correlation of a neuron $l, i$ as $\mathrm{corr}_{l,i} = \max_k \mathrm{corr}_{l,i,k}$. We then identify neurons with top-$K$ absolute values of correlations where $K$ is a hyperparameter.

4. **Interpretation:** We search for sentences in the validation set that maximally or minimally activate the identified neurons. These sentences provide interpretable evidence of potential skills associated with the neuron activations.

We would note that our method is a strict generalization of the method in Wang et al. (2022). In the binary classification setting, if we choose the metric $m$ to map instances of one class to 0 and instances of the other class to 1, then we would select neurons whose activations is predictable of the class labels.

## 3 Experiments

We evaluate our framework on the Qwen 1.5 family of instruction-tuned models (Bai et al., 2023) (with parameter scale 1.8B) and train soft prompts with 20 soft tokens using AdamW optimizer with learning rate 3e-3. Our experiments are designed to address the following research questions:

- **RQ1:** Can we detect skill-related neurons in a natural language generation task when provided with explicit meta labels?

- **RQ2:** Can our framework effectively disentangle and isolate neurons that are specialized for fine-grained linguistic cues or heuristics within a uniform task setting?

- **RQ3:** Can our framework detect skill neurons without relying on explicit meta-labels?

**Skill-Mix (RQ1).** To answer RQ1, we build upon the prompting framework introduced by Yu et al. (2023), which guides LLMs to generate natural language sequences that require specific linguistic

skills. For instance, the logical skill *spatial reasoning* is defined as the capacity to reason about spatial relationships between objects. We modify this framework to generate question-answer pairs that target one of two skills—spatial reasoning or creating metaphor. An example is shown below:

> **Example Data**
>
> Q: Where is the ball if it is to the right of the box and the box is on the table?
> A: The ball is to the right of the table.

We explicitly prompt GPT-4 to produce pairs that leverage one of these two skills and define the metric function $m$ to check whether the selected skill is spatial reasoning. Our analysis reveals a group of neurons with distinct activation patterns corresponding to different skills (Figure 2). This experiment demonstrates that the skill differences within LLMs can be reflected in sparse neurons, and our method can well detect these crucial neurons when meta labels are available.
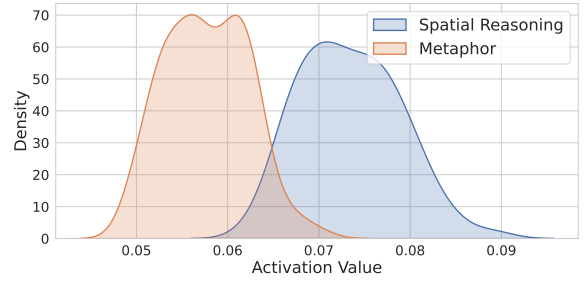


Figure 2: Distribution of activations of the neuron with highest absolute correlation by skill (Skill-Mix). The distribution is interpolated by Kernel density estimation (KDE) based on the empirical distribution of activation over validation set.

**Heuristic Analysis for NLI Systems (HANS) (RQ2).** To investigate RQ2, we apply our framework to the HANS dataset (McCoy et al., 2019), which is specifically designed to reveal whether LLMs rely on simple syntactic shortcuts for natural language inference. Although the overall task remains natural language inference, the dataset contains various fine-grained heuristics. We define our metric $m$ as whether the heuristic is *Lexical Overlap* and probe the neurons accordingly. As shown in Figure 3, the activation distributions of the chosen neurons vary clearly between different heuristics—even for those heuristics that were not directly used during probing. This result confirms

that our framework can identify fine-grained linguistic skill neurons within a single-purpose task.

We further show that high-correlation neurons are sparse on the HANS task. Figure 7 illustrates the distribution of correlation values between each neuron's activation and the HANS heuristic label. While we observe a general clustering of neurons around low correlation values, only a small subset of neurons exhibit substantial correlation, exceeding the threshold marked by the red dashed line (0.43). This sparsity further strengthens the findings from our previous analysis showing that our framework can indeed identify fine-grained linguistic skill neurons even when surrounded by a majority of confounding neurons. We defer other two tasks' correlation distribution to Section A.
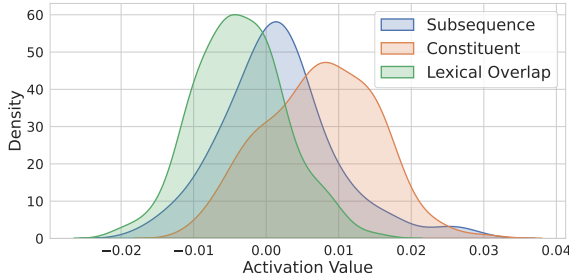


Figure 3: Distribution of activations of the neuron with highest absolute correlation on data with three different heuristics on HANS.
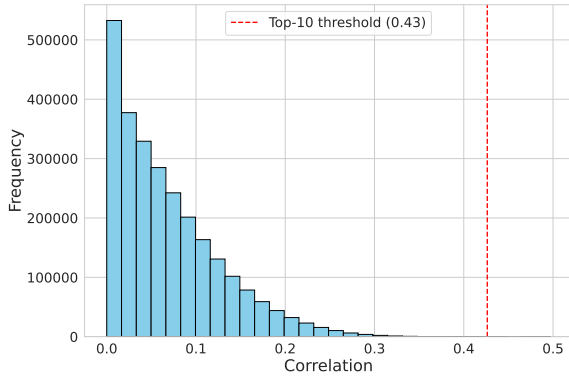


Figure 4: Distribution of correlation values between each neuron's activation and the HANS heuristic label. The red dashed line indicates the top-10 correlation score threshold (0.43).

**Arithmetic Task (RQ3).** For RQ3, we extend our analysis to a task that does not rely on explicit meta labels. We consider a multiple-choice arithmetic problem from BigBench (Srivastava et al., 2023) and define our metric $m$ as the per-sample

loss for each data point. For the selected neuron, we observe that the top 10 sequences with the lowest activations share a common pattern: (1) the question involves multiplication, and (2) the final answer can be determined by considering only the last digit. An example is provided below:

---

**Example Data**

Question: What is 56510 times 52373?
Choices: 16619555, 204563610029, ...,
2959598230
Answer: 2959598230

---

We further verify that the selected neuron corresponds to this subskill by plotting the activation distributions separately for data with and without the identified shortcut (Figure 5). The clear distinction observed between these two groups validates that our framework can detect skill neurons even in tasks lacking explicit meta labels.
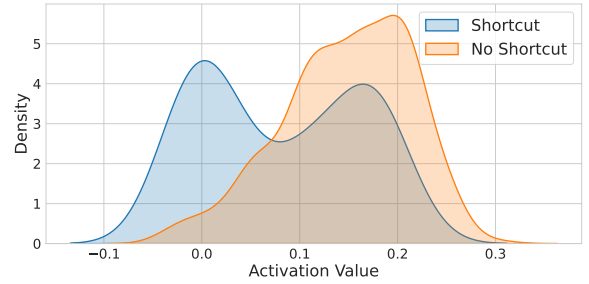


Figure 5: Distribution of activations of the neuron with highest absolute correlation on the data of the Arithmetic task. The shortcut indicates that the correct answer for a multiplication question can be determined solely by the last digit, a pattern automatically discovered by our algorithm.

## 4  Conclusion and Future Work

In this paper, we extend the framework of Wang et al. (2022) on identifying sparse skill-related neurons within LLMs by examining the correlations between neuron activations and a broad range of model behaviors. Experiments on tasks like open-ended question answering, natural language inference, and arithmetic problem-solving demonstrate that our method can consistently identify sparse neurons corresponding to fine-grained linguistic skills and previously unknown heuristics in problem-solving. These findings help understand the skill specialization within LLMs and we encourage future work exploring the causal influence of these identified neurons to model behaviors.

## 5 Limitation

Our method, while effective in uncovering task-specific neuron activations via soft prompt training, depends on the quality of prompt tuning and the availability of clear auxiliary metrics. Moreover, the framework and experiments in this paper are designed to discover the internal activation signatures that are highly *correlational* to model behaviors of interest, and we leave examining the *causality* of these signatures to model behaviors to future work.

## References

Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023. Qwen technical report. *Preprint*, arXiv:2309.16609.

David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. 2017. Network dissection: Quantifying interpretability of deep visual representations. *Preprint*, arXiv:1704.05796.

Leonard Bereska and Efstratios Gavves. 2024. Mechanistic interpretability for ai safety – a review. *Preprint*, arXiv:2404.14082.

Steven Bills, Nick Cammarata, Dan Mossing, Henk Tillman, Leo Gao, Gabriel Goh, Ilya Sutskever, Jan Leike, Jeff Wu, and William Saunders. 2023. Language models can explain neurons in language models. https://openaipublic.blob.core.windows.net/neuron-explainer/paper/index.html.

Dami Choi, Vincent Huang, Kevin Meng, Daniel D Johnson, Jacob Steinhardt, and Sarah Schwettmann. 2024. Scaling automatic neuron description. https://transluce.org/neuron-descriptions.

Wes Gurnee, Neel Nanda, Matthew Pauly, Katherine Harvey, Dmitrii Troitskii, and Dimitris Bertsimas. 2023. Finding neurons in a haystack: Case studies with sparse probing. *Preprint*, arXiv:2305.01610.

Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *Preprint*, arXiv:2104.08691.

Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, Online. Association for Computational Linguistics.

R. Thomas McCoy, Ellie Pavlick, and Tal Linzen. 2019. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. *Preprint*, arXiv:1902.01007.

Alec Radford, Rafal Jozefowicz, and Ilya Sutskever. 2017. Learning to generate reviews and discovering sentiment. *Preprint*, arXiv:1704.01444.

Tilman Räuker, Anson Ho, Stephen Casper, and Dylan Hadfield-Menell. 2023. Toward transparent ai: A survey on interpreting the inner structures of deep neural networks. *Preprint*, arXiv:2207.13243.

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew Dai, Andrew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinion, Cameron Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, César Ferri Ramírez, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites, Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca

Happé, Francois Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovitch-López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin, Hinrich Schütze, Hiromu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B. Simon, James Koppel, James Zheng, James Zou, Jan Kocoń, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John U. Balis, Jonathan Batchelder, Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh D. Dhole, Kevin Gimpel, Kevin Omondi, Kory Mathewson, Kristen Chiafullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros Colón, Luke Metz, Lütfi Kerem Şenel, Maarten Bosma, Maarten Sap, Maartje ter Hoeve, Maheen Farooqi, Manaal Faruqui, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramírez Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin Potthast, Matthew L. Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna Baitemirova, Melody Arnaud, Melvin McElrath, Michael A. Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube, Michał Swędrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Cain, Mimee Xu, Mirac Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal, Moin Aminnaseri, Mor Geva, Mozhdeh Gheini, Mukund Varma T, Nanyun Peng, Nathan A. Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter Chang, Peter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millière, Rhythm Garg, Richard Barnes, Rif A. Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Ronan LeBras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov, Ryan Chi, Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Mohammad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R. Bowman, Samuel S. Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upadhyay, Shyamolima, Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, Stefan Divic, Stefano Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven T. Piantadosi, Stuart M. Shieber, Summer Misherghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsu Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Titus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Ramasesh, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang, and Ziyi Wu. 2023. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Preprint*, arXiv:2206.04615.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. Attention is all you need. *Preprint*, arXiv:1706.03762.

Xiaozhi Wang, Kaiyue Wen, Zhengyan Zhang, Lei Hou, Zhiyuan Liu, and Juanzi Li. 2022. Finding skill neurons in pre-trained transformer-based language models. *Preprint*, arXiv:2211.07349.

Dingli Yu, Simran Kaur, Arushi Gupta, Jonah Brown-Cohen, Anirudh Goyal, and Sanjeev Arora. 2023. Skill-mix: a flexible and expandable family of evaluations for ai models. *Preprint*, arXiv:2310.17567.

# A Additional Experiments

We show the distribution of correlation scores of neurons for SkillMix and Arithmetic tasks here. We observe that the selected neurons' activation has high correlations with the auxiliary metrics and only a very sparse subset of neurons has the same level of correlations.
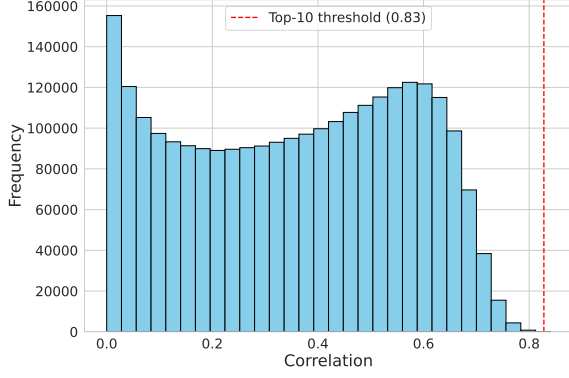


Figure 6: Distribution of correlation values between each neuron's activation and the skill label (for the 1.4B model). The red dashed line indicates the top-10 threshold (0.83).
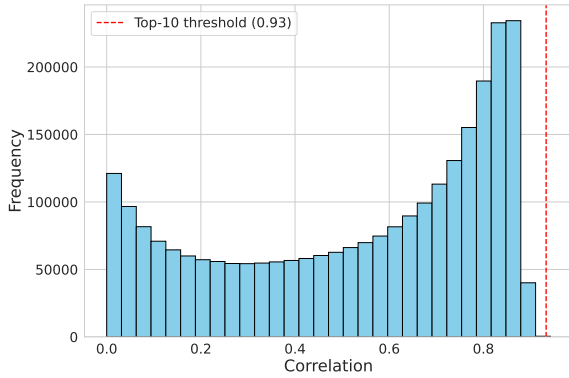


Figure 7: Distribution of correlation values between each neuron's activation and the arithmetic validation sample loss (for the 1.4B model). The red dashed line indicates the top-10 threshold (0.93).