
Universal Neurons in GPT-2: Emergence, Persistence, and Functional Impact

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We investigate the phenomenon of neuron universality in independently trained
2 GPT-2 Small models, examining how these universal neurons—neurons with con-
3 sistently correlated activations across models—emerge and evolve throughout
4 training. By analyzing five GPT-2 models at three checkpoints (100k, 200k, 300k
5 steps), we identify universal neurons through pairwise correlation analysis of acti-
6 vations over a dataset of 5 million tokens. Ablation experiments reveal significant
7 functional impacts of universal neurons on model predictions, measured via loss
8 and KL divergence. Additionally, we quantify neuron persistence, demonstrat-
9 ing high stability of universal neurons across training checkpoints, particularly
10 in deeper layers. These findings suggest stable and universal representational
11 structures emerge during neural network training.

12 1 Introduction

13 Large language models (LLMs) exhibit remarkable generalization but remain difficult to interpret
14 [1]. However, neural networks are fully observable and deterministic, allowing us to record and
15 manipulate internal components such as neuron activations [2]. This presents a rare opportunity to
16 reverse-engineer their internal mechanism. An important open question regarding interpretability is
17 whether models independently trained on the same task converge on similar internal structures—a
18 notion termed the *universality hypothesis* [3]. Universality, if established, offers stable interpretability
19 targets and aids transfer learning.

20 We examine this hypothesis by analyzing five GPT-2 models trained from scratch, tracking when
21 universal neurons—units with highly correlated activations across models [4]—emerge, their stability
22 over training, and their causal role.

23 Contributions:

- 24 • **Emergence Analysis:** We provide the first systematic study of how universal neurons
25 emerge during training, showing that they appear early and grow steadily, especially in
26 deeper layers.
- 27 • **Persistence Quantification:** We quantify the stability of universal neurons across training
28 checkpoints, finding that over 80% remain universal in subsequent stages.
- 29 • **Functional Role via Ablation:** We demonstrate that ablating universal neurons significantly
30 increases loss and KL divergence, confirming their causal importance to model predictions.
- 31 • **Layer-wise Characterization:** We show that first-layer universal neurons disproportionately
32 affect output distributions, suggesting they encode critical low-level information.

2 Related Work

Universality and Cross-Model Consistency. Early studies reported limited direct neuron matching [5]. However, recent work identifies universal neurons with consistent semantic features across independently trained GPT-2 models [4]. These universal neurons are also shown to correspond to interpretable, semantically meaningful features [4]. This provides evidence that some circuits are consistently discovered across training runs, supporting the hypothesis of shared representational scaffolding.

Representation Similarity. Because neurons are often polysemantic [2], direct comparison is difficult. (author?) [6] addressed this by learning sparse features with autoencoders and found significant alignment of feature dimensions across models. Algorithmic behaviors also show cross-architecture consistency, indicating broader universality [7, 8].

Emergence and Stability. Works such as the lottery ticket hypothesis [9] and canonical correlation studies [10] show that networks form persistent representational patterns within the first few training epochs. Theoretical analyses further support the idea that networks rapidly learn dominant features that are gradually refined [11]. These results motivate our focus on the *emergence* and *persistence* of universal neurons throughout training.

3 Method

We analyze five GPT-2 Small models at checkpoints 25%, 50%, and 75% of training (100k, 200k, 300k steps). Neuron activations are extracted over 5M tokens from the Pile dataset [12].

Identifying Universal Neurons via Correlation. Following (author?) [4], we compute Pearson correlations between neurons across model pairs. Let $\mathbf{a}_k^{(m,c)} \in \mathbb{R}^n$ denote the activation vector of neuron k in model m at checkpoint c over n token positions. The Pearson correlation between neurons defined by $\mathbf{a}_k^{(m_1,c)}$ and $\mathbf{a}_\ell^{(m_2,c)}$ is:

$$\rho_{k,\ell}^{(m_1,m_2,c)} = \frac{\mathbb{E} \left[(\mathbf{a}_k^{(m_1,c)} - \mu_k)(\mathbf{a}_\ell^{(m_2,c)} - \mu_\ell) \right]}{\sigma_k \sigma_\ell}$$

where μ_k and σ_k are mean and standard deviations of the activation vector $\mathbf{a}_k^{(m,c)}$ computed across a 5 million token dataset of the uncopyrighted Pile HuggingFace dataset [12]. We compute the excess correlation for a neuron k with respect to a model m_2 at checkpoint c as:

$$\varrho_{k,m_2,c} = \left(\max_{\ell \in N(m_2)} \rho_{k,\ell}^{m_1,m_2,c} - \max_{\ell \in N_R(m_2)} \bar{\rho}_{k,\ell}^{m_1,m_2,c} \right)$$

where $\bar{\rho}_{k,\ell}^{m_1,m_2,c}$ is the Pearson correlation between neuron k in model m_1 and neuron ℓ in a randomly rotated version of the layer from model m_2 , all at a checkpoint c . This rotation is constructed by multiplying the matrix of activations in that layer with a random Gaussian matrix, as described in (author?) [4]. The purpose of this transformation of activation vectors is to eliminate any privileged basis and establish a baseline for comparison [4]. Neurons exceeding an excess correlation of 0.5 for a model j are labeled *universal*. We also adjust this threshold to 0.4 and 0.6 to verify robustness. Tracking ϱ_i across checkpoints and models allows us to observe universality emerge over training.

We used five GPT-2 Small models, models a through e¹. We selected model a as the reference and computed Pearson correlations between its neurons and those in each of the other models (b, c, d, e). This yields four distinct sets of universal neurons.

¹from stanford-crfm at HuggingFace

69 **Persistence Across Training Checkpoints.** We evaluate whether neurons remain universal over
70 time by computing:

$$P_{\text{persist}} = P(\text{univ. at } t_2 \mid \text{univ. at } t_1)$$

71 across training step intervals (e.g., 100k→200k, 200k→300k). To localize this further, we stratify by
72 transformer layer ℓ :

$$P_{\text{persist}}(\ell) = P(\text{univ.}_{t_2} \mid \text{univ.}_{t_1}, \text{layer} = \ell)$$

73 **Ablation Studies and Functional Role.** To test functional significance, we ablate universal and
74 control (non-universal) neurons during inference by zeroing their MLP outputs. We then measure
75 changes in loss and KL divergence of the softmax distributions before and after ablation. We compute
76 KL divergence as

$$\text{KL}(P||Q) = \frac{1}{|T|} \sum_{t \in T} \sum_{x \in V} P_t(x) \log \frac{P_t(x)}{Q_t(x)}$$

77 where T is the set of token positions, V is the output vocabulary (set of all possible tokens), and P, Q
78 are original and ablated softmax distributions of output logits at position t , respectively. We also
79 perform a sensitivity analysis by repeating the experiment with relaxed thresholds for determining
80 universality (e.g., 0.4 and 0.6).

81 4 Results

82 **Emergence of Universal Neurons.** Universal neurons emerge early, increasing consistently through
83 training, notably in deeper layers (Table 1). At early checkpoints (100k steps), fewer than 5% of
84 neurons meet the universality criterion (0.5 threshold), but this fraction increases steadily to nearly
85 6% by 300k steps. Adjusting the universality threshold to 0.4 or 0.6 shows consistent trends.

86 **Persistence of Universal Neurons.** We assess how consistently universal neurons remain universal
87 as training progresses by computing their conditional persistence across three intervals: 100k→200k,
88 200k→300k, and 100k→300k steps.

89 Figure 1 shows layer-wise persistence rates across these intervals. We find that universal neurons
90 are highly stable over time, especially in later layers. Layers 10 and 11 consistently exceed 90%
91 persistence, while mid-layer neurons (e.g., layers 3–5) show greater volatility. Persistence from
92 100k→300k is slightly lower overall than adjacent intervals, reflecting gradual representational drift.

93 These results support the hypothesis that universal neurons, particularly in deeper layers, encode
94 stable and task-relevant features that solidify as training proceeds. For completeness, we include
95 detailed layer-wise persistence plots in Appendix A.1 that further proves the hypothesis.

96 **Ablating Universal Neurons.** To test the functional importance of universal neurons, we ablate
97 them by zeroing their MLP outputs during inference and measure the resulting change in model
98 predictions using KL divergence.

99 Figure 2 shows that ablating all universal neurons (with excess correlation > 0.5) leads to a substantial
100 shift in the output distribution, indicating a significant disruption in the model’s predictions. In

Threshold	100k	200k	300k
0.4	9.58	10.99	11.33
0.5	4.74	5.56	5.71
0.6	2.00	2.35	2.41

Table 1: Percentage of universal neurons at different thresholds and checkpoints, averaged across models.

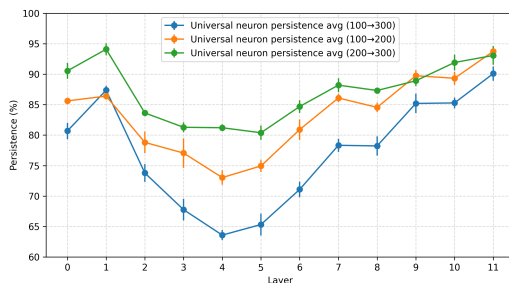


Figure 1: Persistence of universal neurons aggregated by layer (Y-axis begin from 60%).

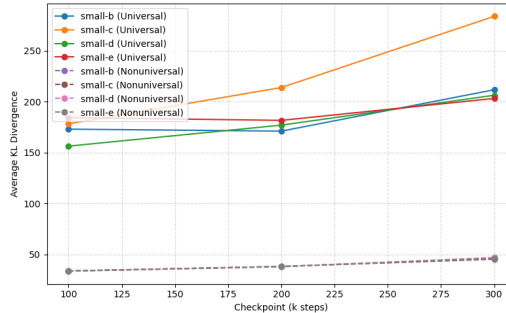


Figure 2: KL divergence of ablated universal (solid) vs. non-universal (dashed) neurons. Colors indicate reference models.

contrast, ablating all non-universal neurons produces minimal change across all checkpoints and models.

These results demonstrate that universal neurons are not only shared across models but also causally important to inference. Considering that only about 5% of all neurons are universal, our results strongly support their role as core components of the model’s learned algorithm.

For completeness, we report supplementary ablation experiments in Appendix B. These include loss change and different thresholds (0.4 and 0.6). In addition to the above findings, layer-wise ablation(Appendix B.2) reveals that ablating universal neurons in the first layer causes a disproportionately large increase in both KL divergence and loss—far exceeding the impact observed in deeper layers. This suggests that early-layer universal neurons play a particularly critical role in shaping the model’s final predictions.

5 Discussion and Conclusion

Findings In this paper, we explore how universal neurons - neurons with high correlations across models - emerge early on in training and persist throughout checkpoints. We found that universal neurons have high functional significance, as ablating them results in higher loss and KL divergence than non-universal neuron ablations. Neurons remain consistent across training checkpoints, with later layers having the higher persistence on average. Trends in universality continue to remain stable despite threshold adjustment (0.4 and 0.6).

Limitations We only studied small models of a few hundred million parameters and monitored activations produced from a data subset of 5 million tokens, which is relatively small. Moreover, we only studied correlations between individual neurons as opposed to families of neurons or higher order circuits, which could offer more interpretable findings.

Future Work To further gain understanding of the impact of universal neurons across families, it would be interesting to examine ablations for families of universal neurons and how the loss varies. A wider selection of experiments could lead to greater insight, for instance, monitoring the effects of activation patching over some training data.

References

- [1] Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, S. Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen A. Creel, Jared Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren E. Gillespie, Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas F. Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, O. Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Kuditipudi,

- 137 Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa
138 Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir P. Mirchandani, Eric
139 Mitchell, Zanele Munyikwa, Suraj Nair, Avani Narayan, Deepak Narayanan, Benjamin
140 Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, J. F. Nyarko, Giray Ogut,
141 Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher
142 Potts, Aditi Raghunathan, Robert Reich, Hongyu Ren, Frieda Rong, Yusuf H. Roohani, Camilo
143 Ruiz, Jack Ryan, Christopher R’e, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy
144 Shih, Krishna Parasuram Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian
145 Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie,
146 Michihiro Yasunaga, Jiaxuan You, Matei A. Zaharia, Michael Zhang, Tianyi Zhang, Xikun
147 Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. On the opportunities and
148 risks of foundation models. *ArXiv*, 2021.
- 149 [2] Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Con-
150 erly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, Robert Lasenby, Yifan Wu,
151 Shauna Kravec, Nicholas Schiefer, Tim Maxwell, Nicholas Joseph, Zac Hatfield-Dodds, Alex
152 Tamkin, Karina Nguyen, Brayden McLean, Josiah E Burke, Tristan Hume, Shan Carter,
153 Tom Henighan, and Christopher Olah. Towards monosemanticity: Decomposing language
154 models with dictionary learning. *Transformer Circuits Thread*, 2023. [https://transformer-](https://transformer-circuits.pub/2023/monosemantic-features/index.html)
155 [circuits.pub/2023/monosemantic-features/index.html](https://transformer-circuits.pub/2023/monosemantic-features/index.html).
- 156 [3] Junxuan Wang, Xuyang Ge, Wentao Shu, Qiong Tang, Yunhua Zhou, Zhengfu He, and Xipeng
157 Qiu. Towards universality: Studying mechanistic similarity across language model architectures.
158 *arXiv preprint arXiv:2410.06672*, Oct 2024.
- 159 [4] Wes Gurnee, Theo Horsley, Zifan Carl Guo, Tara Rezaei Kheirkhah, Qinyi Sun, Will Hathaway,
160 Neel Nanda, and Dimitris Bertsimas. Universal neurons in gpt2 language models. *arXiv preprint*
161 *arXiv:2401.12181*, 2024.
- 162 [5] Simon Kornblith, Mohammad Norouzi, Honglak Lee, and Geoffrey Hinton. Similarity of
163 neural network representations revisited. In Kamalika Chaudhuri and Ruslan Salakhutdinov,
164 editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of
165 *Proceedings of Machine Learning Research*, pages 3519–3529. PMLR, 09–15 Jun 2019.
- 166 [6] Michael Lan, Philip Torr, Austin Meek, Ashkan Khakzar, David Krueger, and Fazl Barez.
167 Sparse autoencoders reveal universal feature spaces across large language models, 2025.
- 168 [7] Jack Lindsey, Wes Gurnee, Emmanuel Ameisen, Brian Chen, Adam Pearce, Nicholas L. Turner,
169 Craig Citro, David Abrahams, Shan Carter, Basil Hosmer, Jonathan Marcus, Michael Sklar,
170 Adly Templeton, Trenton Bricken, Callum McDougall, Hoagy Cunningham, Thomas Henighan,
171 Adam Jermyn, Andy Jones, Andrew Persic, Zhenyi Qi, T. Ben Thompson, Sam Zimmerman,
172 Kelley Rivoire, Thomas Conerly, Chris Olah, and Joshua Batson. On the biology of a large
173 language model. *Transformer Circuits Thread*, 2025.
- 174 [8] Emmanuel Ameisen, Jack Lindsey, Adam Pearce, Wes Gurnee, Nicholas L. Turner, Brian Chen,
175 Craig Citro, David Abrahams, Shan Carter, Basil Hosmer, Jonathan Marcus, Michael Sklar,
176 Adly Templeton, Trenton Bricken, Callum McDougall, Hoagy Cunningham, Thomas Henighan,
177 Adam Jermyn, Andy Jones, Andrew Persic, Zhenyi Qi, T. Ben Thompson, Sam Zimmerman,
178 Kelley Rivoire, Thomas Conerly, Chris Olah, and Joshua Batson. Circuit tracing: Revealing
179 computational graphs in language models. *Transformer Circuits Thread*, 2025.
- 180 [9] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable
181 neural networks. In *International Conference on Learning Representations*, 2019.
- 182 [10] Maithra Raghu, Justin Gilmer, Jason Yosinski, and Jascha Sohl-Dickstein. Svcca: Singular
183 vector canonical correlation analysis for deep learning dynamics and interpretability. In I. Guyon,
184 U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors,
185 *Advances in Neural Information Processing Systems 30*, pages 6076–6085. Curran Associates,
186 Inc., 2017.
- 187 [11] Andrew M. Saxe, James L. McClelland, and Surya Ganguli. A mathematical theory of semantic
188 development in deep neural networks. *Proceedings of the National Academy of Sciences*,
189 116(23):11537–11546, 2019.

190 [12] monology. Pile uncopyrighted: A cleaned version of the pile dataset with copyrighted subsets
 191 removed. <https://huggingface.co/datasets/monology/pile-uncopyrighted>, 2023.
 192 Accessed: 2025-06-26.

193 A Persistence of Universal Neurons over Checkpoints

194 A.1 Global Universal Neuron Persistence

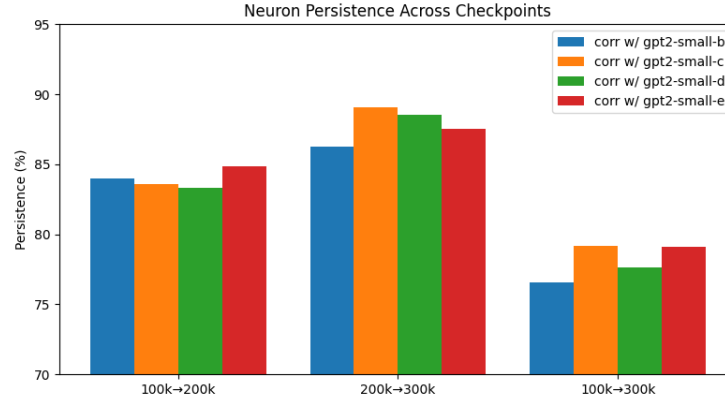


Figure 3: Layer-aggregated persistence of universal neurons across training checkpoints for each model. Bars represent the percentage of universal neurons at the earlier checkpoint that remain universal at the later one. Later-stage intervals (e.g., 200k→300k) exhibit higher persistence, indicating stabilization of universal features over training.

195 A.2 Detailed Layer-wise Universal Neuron Persistence

196 All models exhibit a U-shaped trend: lower persistence in middle layers (2–5) and higher stability
 197 in both early and especially late layers. This suggests that early and late layers encode more stable, model-aligned features during training.

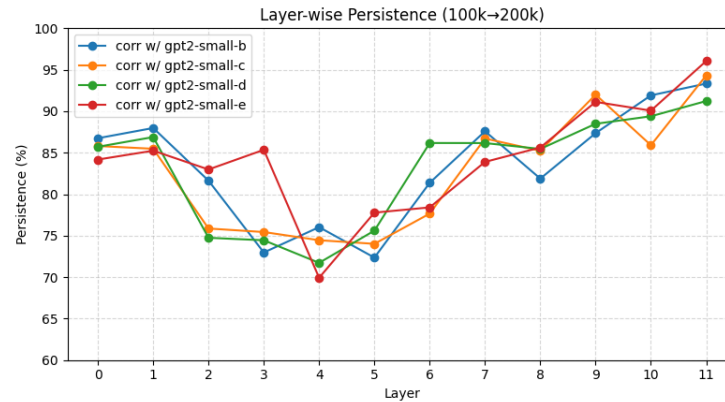


Figure 4: Layer-wise persistence of universal neurons from checkpoint 100k to 200k.

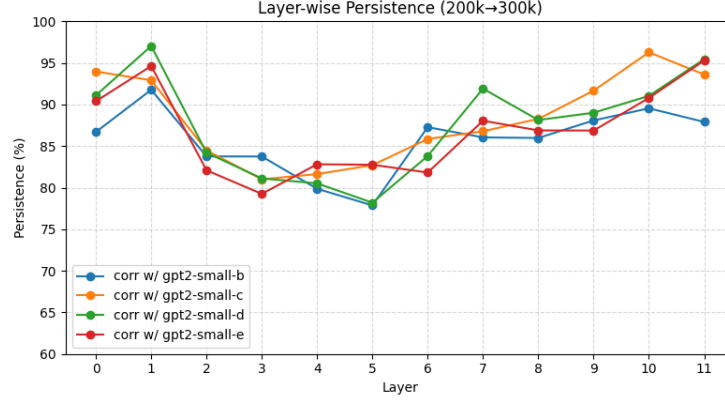


Figure 5: Layer-wise persistence of universal neurons from checkpoint 200k to 300k.

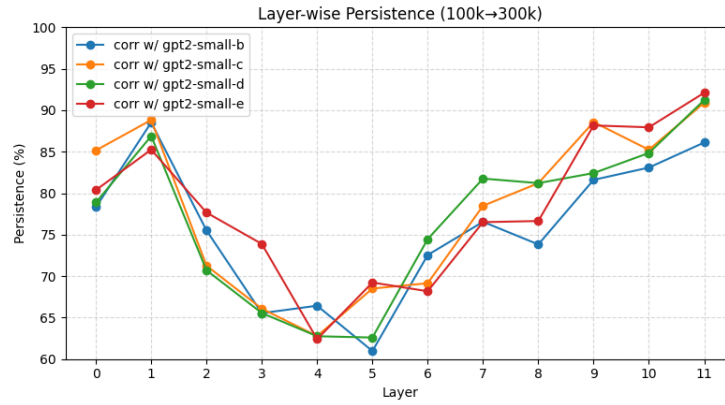


Figure 6: Layer-wise persistence of universal neurons from checkpoint 100k to 300k.

199 B Ablation Experiment Results

200 B.1 Loss Increase From Ablating All Universal/Non-universal Neurons

201 The increase in loss is computed as the difference in model loss before and after neuron ablation.
 202 Across all training checkpoints, ablating universal neurons results in a substantially greater increase
 203 in loss compared to non-universal neurons. The accompanying bar graphs support the primary claim
 204 presented in Section 4.

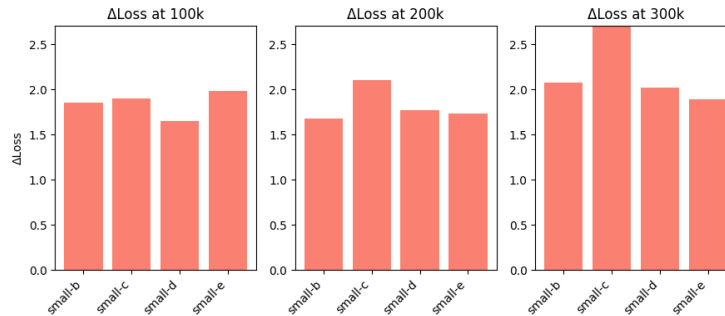


Figure 7: Loss increase from ablating universal neurons over training steps.

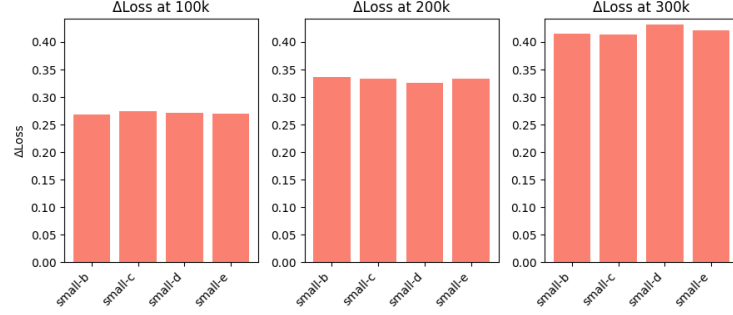


Figure 8: Loss increase from ablating non-universal neurons over training steps.

205 B.2 Comparative Layer-Wise Effects of Universal vs. Non-Universal Neuron Ablation

206 At each training checkpoint, we ablate all universal or non-universal neurons within a specific layer
 207 and evaluate the resulting change in model output using KL divergence and loss increase. Figures
 208 below present the average KL divergence and loss difference (ablated minus original) across five
 209 models for both universal and non-universal neuron ablations.

210 Ablating universal neurons consistently leads to a greater increase in KL divergence and loss compared
 211 to non-universal neurons, indicating their stronger causal role in shaping model predictions. Notably,
 212 the first layer shows the most pronounced sensitivity to ablation, suggesting that early-layer universal
 213 neurons encode particularly critical information. In contrast, non-universal neuron ablation results in
 minor and often negligible effects across layers and checkpoints.

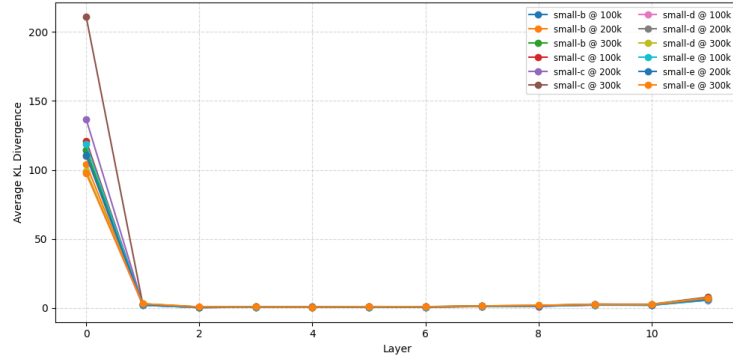


Figure 9: Layer-wise KL divergence from ablating universal neurons.

214

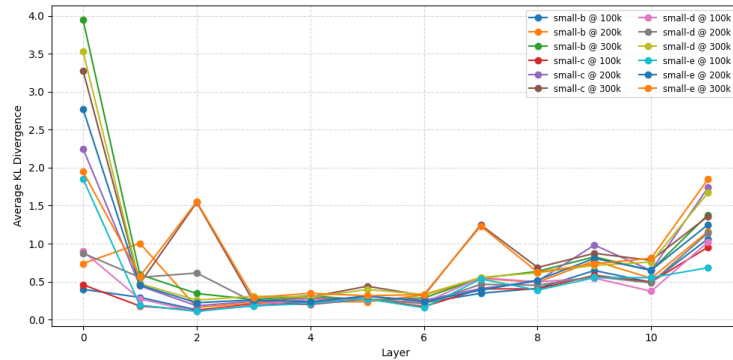


Figure 10: Layer-wise KL divergence from ablating non-universal neurons.

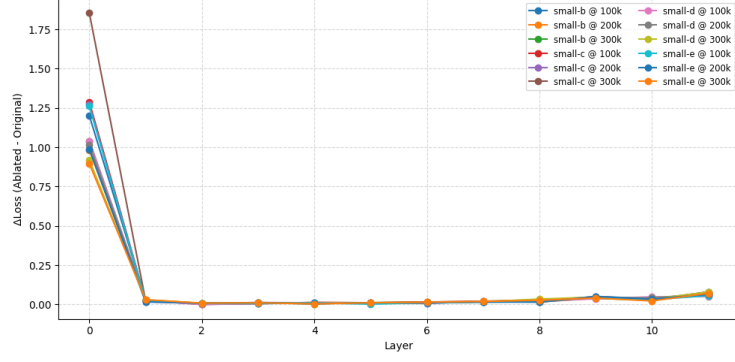


Figure 11: Layer-wise loss increase from ablating universal neurons.

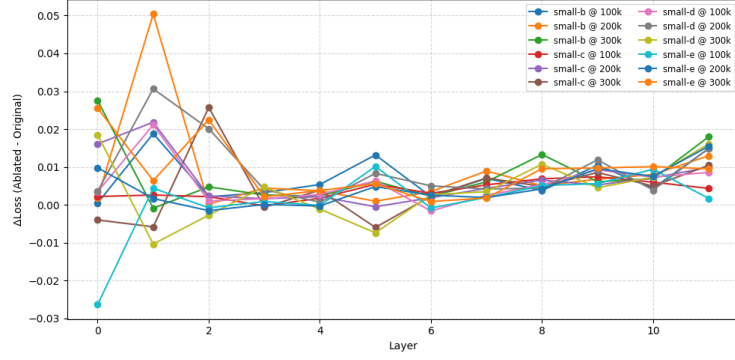


Figure 12: Layer-wise loss increase from ablating non-universal neurons.

215 C Robustness Analysis Under Varying Universality Thresholds(Excess 216 Correlation of 0.4 and 0.6)

217 Although the absolute values vary with different universality thresholds (0.4 and 0.6), the over-
218 all trends remain consistent. These results support our primary claims regarding the functional
219 importance of universal neurons.

220 C.1 All Universal/Non-universal Neurons Ablation with 0.4 Thresholding

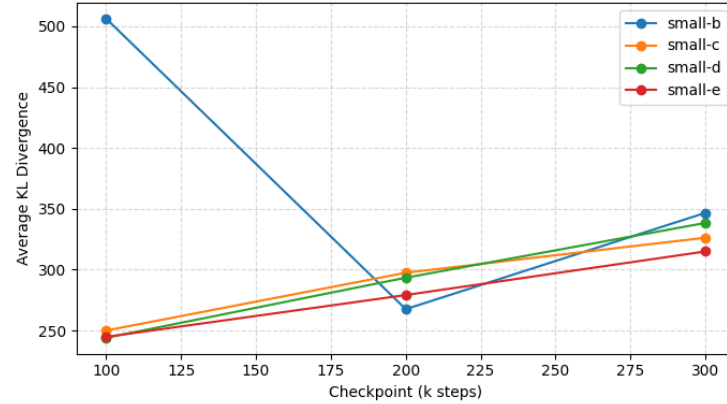


Figure 13: KL divergence from ablating universal neurons over training steps.

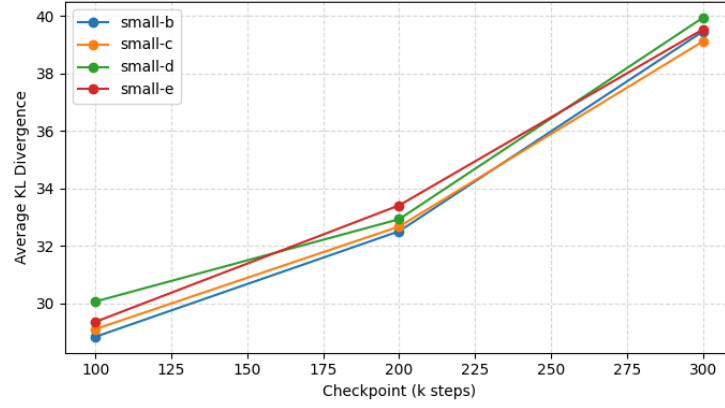


Figure 14: KL divergence from ablating non-universal neurons over training steps.

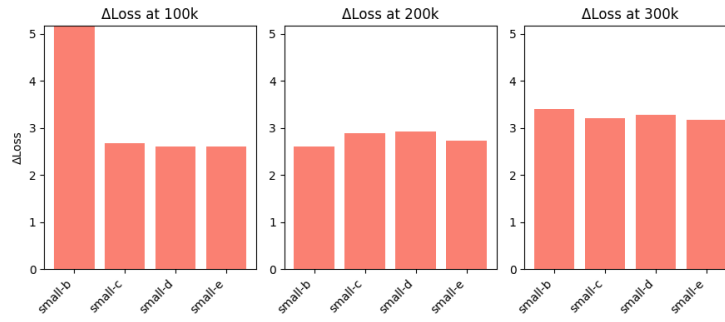


Figure 15: Loss increase from ablating universal neurons over training steps.

221 C.2 All Universal/Non-universal Neurons Ablation with 0.6 Thresholding

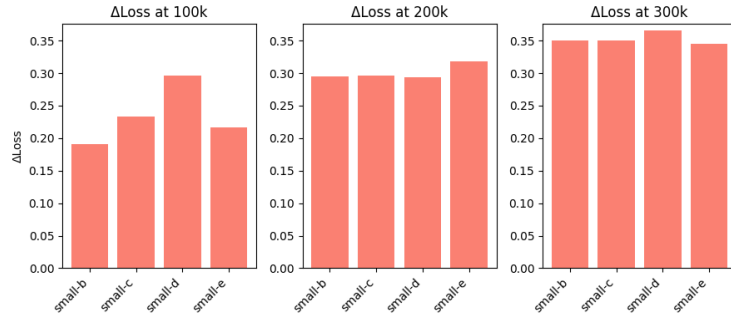


Figure 16: Loss increase from ablating non-universal neurons over training steps.

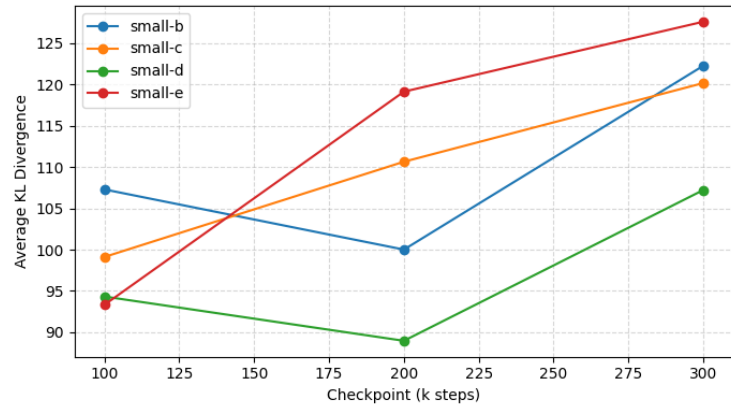


Figure 17: KL divergence from ablating universal neurons over training steps.

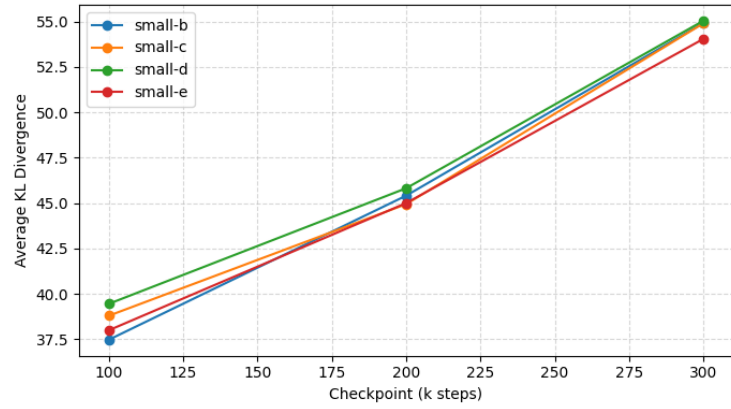


Figure 18: KL divergence from ablating non-universal neurons over training steps.

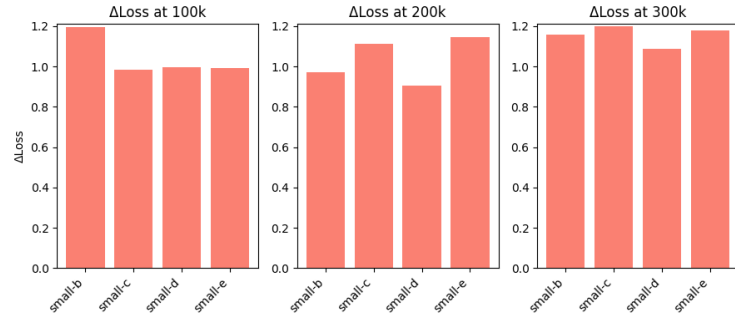


Figure 19: Loss increase from ablating universal neurons over training steps.

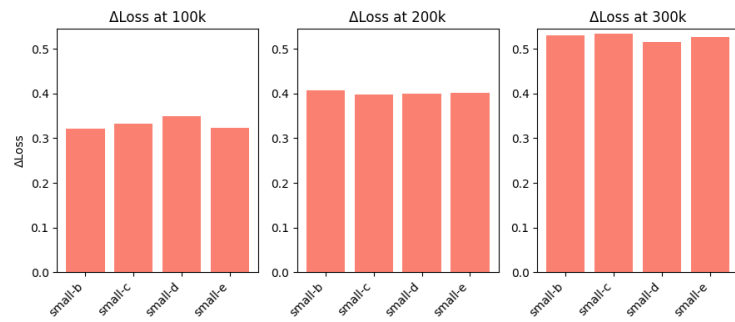


Figure 20: Loss increase from ablating non-universal neurons over training steps.