

---

# No Clustering, No Routing: How Transformers Actually Process Rare Tokens

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Large language models struggle with rare token prediction, yet the mechanisms driv-  
2 ing their specialization remain unclear. Prior work identified specialized “plateau”  
3 neurons for rare tokens following distinctive three-regime influence patterns [7],  
4 but their functional organization is unknown. We investigate this through neuron  
5 influence analyses, graph-based clustering, and attention head ablations in GPT-2  
6 XL and Pythia models. Our findings show that: (1) rare token processing requires  
7 additional plateau neurons beyond the power-law regime sufficient for common  
8 tokens, forming dual computational regimes; (2) plateau neurons are spatially dis-  
9 tributed rather than forming modular clusters; and (3) attention mechanisms exhibit  
10 no preferential routing to specialists. These results demonstrate that rare token  
11 specialization arises through distributed, training-driven differentiation rather than  
12 architectural modularity, preserving context-sensitive flexibility while achieving  
13 adaptive capacity allocation.

## 14 1 Introduction

15 Large language models (LLMs) achieve remarkable performance across diverse tasks, yet their inter-  
16 nal mechanisms for processing different token types remain poorly understood. Rare tokens—those  
17 appearing infrequently in training data pose particular challenges: they often encode critical semantic  
18 information while exhibiting systematically lower prediction accuracy [5]. Prior studies show that  
19 model performance on factual tasks correlates with entity frequency [1] and that truncating rare  
20 tokens in training can lead to model collapse [12], underscoring the importance of understanding how  
21 LLMs process rare events.

22 Recent advances in mechanistic interpretability have revealed specialized circuits for specific lin-  
23 guistic computations, including feed-forward layers acting as key-value memories [4] and sparse,  
24 monosemantic representations [3]. In particular, Liu et al. [7] found that rare token processing  
25 exhibits a distinctive three-regime influence distribution: a plateau of highly influential neurons,  
26 followed by a power-law decay and a rapid decay and they argue that this might be related with  
27 neuron coordinated in both activation and weight space.

28 However, while specialized plateau neurons exist for rare tokens, **it is unclear how transformers**  
29 **organize and access these computational resources**. This organizational question connects to  
30 complementary learning systems theory [9, 6], where two frameworks make opposing predictions.  
31 The **modular hypothesis** suggests specialization requires discrete neuron clusters and selective  
32 routing [11], while the **distributed hypothesis** proposes that specialization emerges from parameter-  
33 level differentiation within shared substrates [8].

34 To test these hypotheses, we investigate three questions: (1) whether rare and common tokens  
35 require different computational regimes, (2) whether plateau neurons form spatial clusters, and (3)

whether attention mechanisms selectively route rare tokens to specialists. Our results support the distributed hypothesis: plateau neurons provide additional processing capacity while remaining spatially distributed and accessed through universal attention patterns. This challenges architectural modularity assumptions and demonstrates that sophisticated computational organization can emerge from simple training dynamics.

## 2 Methods

We investigated rare token specialization in GPT-2 XL and Pythia models through three complementary analyses: neuron influence, spatial organization, and attention routing. All analyses focus on the final MLP layer, where prior work observed the characteristic three-regime influence pattern [7].

**Dataset and Token Selection** We sampled 25,088 tokens from the C4 corpus [10]. Tokens were split into rare and common groups based on frequency, using the 50th percentile as a threshold to ensure balanced sample sizes while capturing long-tail effects.

**Neuron Influence Analysis** To quantify individual neuron contributions, we performed mean ablation experiments. Influence of neuron  $n$  was measured as the absolute change in loss after ablation:

$$\text{Influence}(n) = |\mathcal{L}_{\text{ablated}}(n) - \mathcal{L}_{\text{baseline}}|. \quad (1)$$

Neurons were ranked by influence, and a power-law curve was fitted to the distribution. Neurons whose influence significantly exceeded the fit were classified as the plateau regime, while mid- and low-influence neurons followed the standard decay. Comparing rare and common tokens allowed us to assess whether distinct computational regimes emerge.

**Spatial Organization Analysis** To test whether plateau neurons form modular clusters, we constructed correlation networks from neuron activations and applied Louvain community detection [2]. Signed modularity  $Q$  quantifies clustering:

$$Q = \frac{1}{2m} \sum_{ij} \left[ A_{ij} - \frac{k_i k_j}{2m} \right] \delta(c_i, c_j), \quad (2)$$

where  $A_{ij}$  is the correlation-based edge weight,  $k_i$  the node degree, and  $\delta(c_i, c_j)$  indicates same-community membership. Plateau neuron modularity was compared to random baselines to evaluate significance.

**Attention Routing Analysis** To examine whether rare tokens rely on specialized attention routing, we analyzed attention patterns in layers preceding the final MLP. Attention concentration was quantified via Gini coefficients, and correlations between rare and common token attention distributions were computed. Head-specific contributions were measured using ablation:

$$\text{Impact}(h) = \frac{|\text{Activation}_{\text{baseline}} - \text{Activation}_{\text{ablated } h}|}{\text{Activation}_{\text{baseline}}}. \quad (3)$$

Minimal impact from single-head ablations, compared to large drops from full-layer ablation, indicates that plateau neurons integrate signals from multiple heads rather than being selectively targeted. All analyses included statistical controls and significance testing to ensure that observed patterns reflect genuine specialization rather than noise or sampling variability.

## 3 Results

**Rare and Common Tokens Show Distinct Influence Patterns** Figure 1 shows that rare and common tokens engage distinct computational regimes. For common tokens, neuron influence follows a well-fitted power-law:

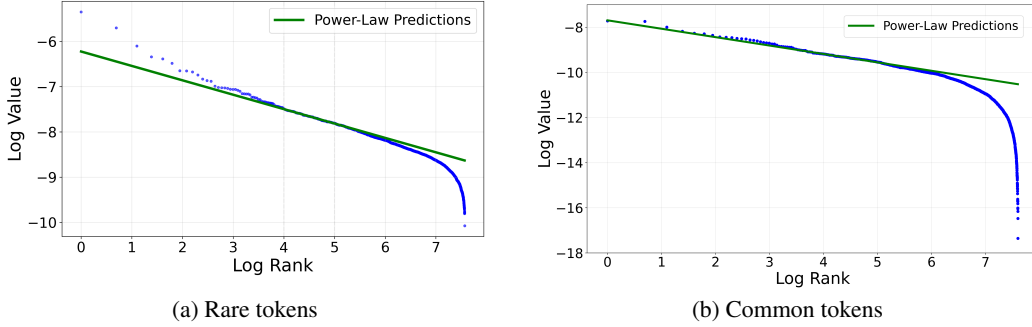


Figure 1: Neuron influence distributions for rare vs. common tokens. Rare tokens exhibit a plateau of specialist neurons, a power-law region, and a rapid decay region. Common tokens follow a pure power-law distribution without a plateau regime.

$$\log |\Delta\mathcal{L}| \approx -\kappa \log(\text{rank}) + \beta, \quad (4)$$

with  $\kappa = 1.84 \pm 0.12$  ( $R^2 = 0.94$ ). Deviations from this fit remain small ( $|\delta| < 0.1$  for 95% of neurons), confirming that common tokens primarily rely on distributed scaling.

In contrast, rare tokens systematically deviate from power-law behavior. Among the top 15–20 neurons, we observe a clear plateau regime with positive deviations ( $\delta > 0.5$ ) relative to the fitted curve. Beyond the plateau, mid-ranked neurons follow a similar power-law decay, while low-influence neurons show rapid signal attenuation. This establishes a dual-regime structure for rare tokens, consistent across both GPT-2 XL and Pythia families.

**Plateau Neurons Are Spatially Distributed** We next tested whether plateau neurons form spatial clusters or are distributed across the MLP layer. Using graph-based Louvain community detection, we computed signed modularity scores for excitatory and inhibitory plateau neurons and compared them to random controls (Table 1).

Across both models, all plateau neuron groups had modularity scores comparable to random baselines (Pythia-410M: 0.03–0.05 vs. control 0.04; GPT-2 XL: 0.07–0.11 vs. control 0.09), with no statistically significant clustering (p-values 0.42–0.84). Spectral clustering produced consistent null results ( $Q_{\text{plateau}} = 0.08 \pm 0.05$  vs.  $Q_{\text{control}} = 0.07 \pm 0.06$ ,  $p = 0.41$ , Mann–Whitney  $U$  test). Together, these analyses indicate that plateau neurons are spatially distributed rather than forming discrete modules.

Table 1: Community detection results show no significant clustering of plateau neurons compared to random baselines across model scales.

| Model       | Neuron Group         | Signed Modularity | Communities   | p-value |
|-------------|----------------------|-------------------|---------------|---------|
| Pythia-410M | Plateau (Excitatory) | $0.05 \pm 0.04$   | $1.8 \pm 0.5$ | 0.67    |
|             | Plateau (Inhibitory) | $0.03 \pm 0.05$   | $1.9 \pm 0.6$ | 0.84    |
|             | Random Control       | $0.04 \pm 0.04$   | $1.9 \pm 0.4$ | –       |
| GPT-2 XL    | Plateau (Excitatory) | $0.11 \pm 0.06$   | $2.3 \pm 0.7$ | 0.42    |
|             | Plateau (Inhibitory) | $0.07 \pm 0.05$   | $2.1 \pm 0.5$ | 0.73    |
|             | Random Control       | $0.09 \pm 0.05$   | $2.2 \pm 0.6$ | –       |

**No Evidence for Selective Attention Routing** To examine whether rare tokens rely on specialized attention routing, we analyzed attention distributions in layers preceding the final MLP. Correlations between rare and common tokens were high ( $r = 0.89 \pm 0.07$ ), and attention concentration (Gini coefficient) did not differ significantly ( $G_{\text{rare}} = 0.34 \pm 0.05$  vs.  $G_{\text{common}} = 0.32 \pm 0.04$ ,  $p = 0.43$ , t-test).

Systematic head ablation experiments further confirmed distributed access (Table 2). Single-head ablations caused small, statistically similar reductions in plateau activation (effect sizes 0.26–0.31),

97 whereas removing all heads in a layer produced large drops (-42% to -45%). These results indicate  
 98 that plateau neurons integrate signals from multiple heads rather than relying on dedicated routing  
 99 circuits.

100 Overall, our results show that rare tokens (1) recruit plateau neurons not activated by common tokens,  
 101 establishing distinct computational regimes; (2) engage neurons that are spatially distributed rather  
 102 than clustered; and (3) access these neurons via distributed attention mechanisms with no selective  
 103 routing. This supports a general principle of distributed specialization in transformers rather than  
 104 modular organization.

Table 2: Attention head ablation shows distributed dependency patterns across model scales. Individual heads show minimal impact on plateau neuron activation.

| Model       | Ablation Target          | Plateau Activation Change | Effect Size | p-value |
|-------------|--------------------------|---------------------------|-------------|---------|
| Pythia-410M | Single Head (max impact) | $-7.8\% \pm 2.3\%$        | 0.29        | 0.03    |
|             | Random Head (baseline)   | $-7.1\% \pm 2.9\%$        | 0.26        | 0.04    |
|             | All Heads                | $-42.1\% \pm 5.2\%$       | 1.74        | <0.001  |
|             | Control (non-attention)  | $-1.3\% \pm 1.6\%$        | 0.05        | 0.71    |
| GPT-2-XL    | Single Head (max impact) | $-8.2\% \pm 2.1\%$        | 0.31        | 0.02    |
|             | Random Head (baseline)   | $-7.4\% \pm 3.2\%$        | 0.28        | 0.03    |
|             | All Heads                | $-45.3\% \pm 4.7\%$       | 1.82        | <0.001  |
|             | Control (non-attention)  | $-1.1\% \pm 1.8\%$        | 0.04        | 0.67    |

## 105 4 Discussion and Limitations

106 Our results demonstrate that transformers handle rare tokens through distributed, training-driven  
 107 specialization rather than modular architectural design. Rare tokens recruit additional high-influence  
 108 plateau neurons that are largely inactive during common token processing, establishing dual com-  
 109 putational regimes. These plateau neurons are spatially distributed across the MLP layer and are  
 110 accessed through the same attention patterns used by all tokens, with no evidence of selective routing.  
 111 Together, these findings indicate that transformers achieve rare token specialization by differenti-  
 112 ating parameters within shared computational substrates, preserving flexible and context-sensitive  
 113 processing without requiring discrete modules.

114 This distributed organization provides insight into the scalability and robustness of transformer  
 115 architectures. By leveraging parameter-level specialization instead of hardwired modular structures,  
 116 transformers can allocate computational capacity adaptively, avoiding bottlenecks while handling  
 117 low-frequency but semantically critical tokens. Our analysis suggests that mixture-of-experts style  
 118 routing may be unnecessary for rare token processing, as universal connectivity already enables  
 119 effective integration of specialized neurons.

120 Despite these insights, our study has several limitations. First, our findings are correlational: we  
 121 observe robust patterns of plateau neuron recruitment and distributed access, but we do not establish  
 122 causal mechanisms or developmental trajectories during training. Second, our analysis focuses  
 123 on GPT-2 and Pythia families and specific network components, namely the final MLP layer and  
 124 attention heads near the output. Whether these distributed mechanisms generalize across architectures,  
 125 model scales, or modalities remains an open question. Third, ablation methods measure necessity but  
 126 may overlook distributed redundancy, and token-matching controls cannot fully capture semantic  
 127 complexity.

128 Future work should examine how plateau neurons emerge during training, extend analyses to diverse  
 129 architectures, and employ causal interventions to test their functional necessity. By revealing how rare  
 130 token specialization arises from distributed differentiation rather than modularity, our findings provide  
 131 a foundation for understanding transformer interpretability and for guiding principled architectural  
 132 design.

## References

- [1] E. Akyürek, T. Schick, K. Kawaguchi, M. Antoniak, R. Chen, T. Wang, et al. Towards tracing factual knowledge in language models back to the training data. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2429–2446, 2022.
- [2] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre. Fast unfolding of communities in large networks. *Journal of statistical mechanics: theory and experiment*, 2008(10):P10008, 2008.
- [3] T. Bricken, C. Templeton, and J. Steinhardt. Monosemanticity: Localized features in neural networks and brains. *arXiv preprint arXiv:2310.10999*, 2023.
- [4] M. Geva, R. Schuster, J. Berant, and O. Levy. Transformer feed-forward layers are key-value memories. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5484–5495, 2021.
- [5] N. Kandpal, H. Deng, A. Roberts, E. Wallace, and C. Raffel. Large language models struggle to learn long-tail knowledge. In *International Conference on Machine Learning*, pages 15696–15707. PMLR, 2023.
- [6] D. Kumaran, D. Hassabis, and J. L. McClelland. What learning systems do intelligent agents need? complementary learning systems theory updated. *Trends in Cognitive Sciences*, 20(7): 512–534, 2016.
- [7] J. Liu, H. Wang, and Y. Li. Emergent specialization: Rare token neurons in language models. *arXiv preprint arXiv:2505.12822*, 2025.
- [8] J. L. McClelland, D. E. Rumelhart, P. R. Group, et al. Parallel distributed processing: Explorations in the microstructure of cognition. *Volume 1: Foundations*, 1986.
- [9] J. L. McClelland, B. L. McNaughton, and R. C. O’Reilly. Why there are complementary learning systems in the hippocampus and neocortex: insights from the successes and failures of connectionist models of learning and memory. *Psychological Review*, 102(3):419, 1995.
- [10] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- [11] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *International Conference on Learning Representations*, 2017.
- [12] I. Shumailov, Z. Shumaylov, Y. Zhao, Y. Gal, N. Papernot, and R. Anderson. The curse of recursion: Training on generated data makes models forget. *arXiv preprint arXiv:2305.17493*, 2023.

## 167 A Detailed Methodological Procedures

### 168 A.1 Graph Construction for Spatial Organization Analysis

169 We construct weighted, undirected graphs from neural activation data to assess functional connectivity  
 170 patterns. For a given set of neurons, we compute edge weights using the Pearson correlation coefficient  
 171 between activation vectors of neurons  $i$  and  $j$  across our corpus of 1,000 contexts:

$$w_{ij} = \text{corr}(\mathbf{a}_i, \mathbf{a}_j) \quad (5)$$

172 where  $\mathbf{a}_i$  represents the activation vector for neuron  $i$  across all contexts. The resulting graphs are  
 173 signed, capturing both positive and negative correlations. We apply a threshold of  $|w_{ij}| > 0.1$  to  
 174 remove weak connections and focus on meaningful functional relationships.

### 175 A.2 Community Detection Algorithms

176 **Louvain Algorithm:** Our primary community detection method uses the Louvain algorithm [2],  
 177 which optimizes modularity through iterative local optimization and community aggregation. We run  
 178 the algorithm 100 times with different random seeds and select the partition with highest modularity  
 179 score to ensure robustness.

180 **Spectral Clustering:** As a validation approach, we employ spectral clustering on the graph’s  
 181 normalized Laplacian matrix. We perform eigendecomposition and embed nodes in a low-dimensional  
 182 space ( $k=2$  to  $k=8$  communities) where clusters are identified via k-means clustering. This method is  
 183 particularly effective for detecting complex community structures in signed graphs.

### 184 A.3 Modularity Measures

185 We employ both standard and signed modularity measures to quantify clustering quality:

186 **Standard Modularity:**

$$Q = \frac{1}{2m} \sum_{ij} \left[ A_{ij} - \frac{k_i k_j}{2m} \right] \delta(c_i, c_j) \quad (6)$$

187 **Signed Modularity:** For signed networks, we use the extension:

$$Q_{\text{signed}} = \frac{1}{2m^+} \sum_{ij} \left[ A_{ij}^+ - \frac{k_i^+ k_j^+}{2m^+} \right] \delta(c_i, c_j) - \frac{1}{2m^-} \sum_{ij} \left[ A_{ij}^- - \frac{k_i^- k_j^-}{2m^-} \right] \delta(c_i, c_j) \quad (7)$$

188 where  $A^+$  and  $A^-$  represent positive and negative edge subgraphs respectively.

### 189 A.4 Statistical Validation Procedures

190 **Control Group Generation:** For each experimental group (plateau neurons), we generate 100  
 191 random control groups of the same size from the full neuron population. This controls for baseline  
 192 connectivity patterns and group size effects.

193 **Significance Testing:** We use Mann-Whitney U tests to compare modularity distributions between  
 194 experimental and control groups, with Bonferroni correction for multiple comparisons. Effect sizes  
 195 are calculated using Cohen’s d.

196 **Bootstrap Confidence Intervals:** Modularity scores are estimated with 95

### 197 A.5 Attention Head Ablation Implementation

198 **Ablation Procedure:** For each attention head  $h$  in layers 20-30, we zero the attention weight matrix  
 199  $\mathbf{W}_h$  and forward-propagate to measure downstream effects on plateau neuron activations. This is  
 200 implemented through:

$$\text{Attention}_{\text{ablated}} = \text{softmax} \left( \frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}} \odot \mathbf{M}_h \right) \mathbf{V} \quad (8)$$

201 where  $\mathbf{M}_h$  is a binary mask with zeros for the ablated head.

202 **Control Ablations:** We perform control ablations on randomly selected heads and non-attention  
 203 components (layer norms, feedforward weights) to establish baseline impact levels and ensure  
 204 observed effects are attention-specific.

205 **Activation Patching:** For validation, we implement clean/corrupted activation patching [?] where we  
 206 replace attention outputs with activations from different input contexts to isolate causal contributions.

## 207 A.6 Token Selection and Matching Criteria

208 **Frequency Thresholds:** Rare tokens defined as appearing  $< 100$  times in OpenWebText training  
 209 corpus; common tokens appearing  $> 10,000$  times, based on GPT-2 tokenizer frequency distributions.

210 **Matching Procedure:** Each rare token is matched to a common token controlling for: (1) token length  
 211 ( $\pm 1$  character), (2) part-of-speech tag (using spaCy), (3) sentence position (beginning/middle/end),  
 212 and (4) syntactic role when possible. Matching accuracy verified through linguistic feature analysis.

213 **Context Selection:** For each token pair, we sample 20 diverse contexts ensuring grammatical validity  
 214 and semantic coherence. Contexts are drawn from different domains (news, literature, technical  
 215 writing) to test generalization across linguistic environments.