

# The Achilles' Heel of LLMs: How Altering a Handful of Neurons Can Cripple Language Abilities

Zixuan Qin<sup>1</sup>   Qingchen Yu<sup>2,3</sup>   Kunlin Lyu<sup>1</sup>   Zhaoxin Fan<sup>2,3\*</sup>  
Yifan Sun<sup>1\*</sup>

<sup>1</sup>School of Statistics, Renmin University of China

<sup>2</sup>Beijing Advanced Innovation Center for Future Blockchain and Privacy Computing

<sup>3</sup>School of Artificial Intelligence, Beihang University



## From Human Brains to Large Language Models (LLMs)

- **Neuroscience insight:** Human cognitive functions heavily rely on a small subset of "bottleneck neurons"; damage to these causes immediate capability loss.
- **The fundamental question:** Despite their massive, distributed parameter scale (tens of billions), do LLMs also possess a similarly sparse set of *critical artificial neurons* indispensable for their capabilities?
- **Previous gap:** Most research focuses on weight/activation outliers, lacking neuron-level causal verification of capability collapse.

# Methodology: Perturbation-based Causal Identification

## Stage 1: Neuron Importance Evaluation

Inject controlled Gaussian noise into the input embedding. Generating a ranked list of "potential critical neurons."

## Stage 2: Critical Neuron Identification (Greedy Causal Verification)

Sequentially mask the top-ranked neurons. Monitor for a catastrophic surge in perplexity.

### Algorithm 1 Stage 1: Neuron Importance Evaluation

**Require:** Model  $f_\theta$ , input text  $p$ , noise scale  $\alpha$ , samples  $K$

- 1:  $\mathbf{x} \leftarrow \text{Context}(p)$
- 2:  $A^{\text{clean}} \leftarrow f_\theta(\mathbf{x})$
- 3: **for**  $i = 1$  **to**  $K$  **do**
- 4:  $\tilde{\mathbf{x}}_i \leftarrow \mathbf{x} + \alpha \cdot \epsilon_i$ , where  $\epsilon_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5:  $A_i^{\text{noisy}} \leftarrow f_\theta(\tilde{\mathbf{x}}_i)$
- 6: **end for**
- 7: **for each** neuron  $s$  **do**
- 8:  $\text{Imp}(s) \leftarrow \frac{1}{K} \sum_{i=1}^K |A_s^{\text{clean}} - A_{i,s}^{\text{noisy}}|$
- 9: **end for**
- 10:  $S \leftarrow \text{SortByDescending}(\{s_1, s_2, \dots, s_N\}, \text{Imp})$
- 11: **return** Importance-ranked neuron list  $S$

### Algorithm 2 Stage 2: Critical Neuron Identification

**Require:** Model  $f_\theta$ , sorted neurons  $S$ , input text  $p$ , threshold  $\epsilon$ .

- 1:  $\text{PPL}_{\text{original}} \leftarrow \text{Perplexity}(f_\theta, p)$
- 2:  $M \leftarrow \emptyset$
- 3: **for**  $n = \Delta n, 2\Delta n, \dots, |S|$  **do**
- 4:  $M \leftarrow \text{Top-}n(S)$
- 5: Apply masking:  $\tilde{n}_l^{(i)}(\mathbf{x}) = 0$  for all  $(l, i) \in M$
- 6:  $\text{PPL}_{\text{masked}} \leftarrow \text{Perplexity}(f_\theta^{-M}, p)$
- 7:  $\Delta = \log_{10}(\text{PPL}_{\text{masked}}) - \log_{10}(\text{PPL}_{\text{original}})$
- 8: **if**  $\Delta \geq \epsilon$  **then**
- 9: **break**
- 10: **end if**
- 11: **end for**
- 12: **return** Critical neuron set  $M$

# Finding 1: Ultra-Sparse Vulnerability

**Catastrophic Collapse:** Masking merely **3–10 neurons** ( $\sim 10^{-8}$  of total parameters) completely destroys language ability.

Model Family	Model	Critical Neurons		WikiText-103	C4
		Number	Rate ( $10^{-8}$ )	Original $\rightarrow$ Masked	Original $\rightarrow$ Masked
Llama-3	Llama-3.2-1B-Instruct	4	6.43	17.74 $\rightarrow$ <b><math>1.58 \times 10^6</math></b>	21.15 $\rightarrow$ <b><math>1.41 \times 10^6</math></b>
	Llama-3.2-3B-Instruct	4	3.19	13.36 $\rightarrow$ <b><math>3.79 \times 10^5</math></b>	16.47 $\rightarrow$ <b><math>4.69 \times 10^5</math></b>
	Llama-3-8B-Instruct	5	2.21	10.88 $\rightarrow$ <b><math>1.48 \times 10^6</math></b>	14.04 $\rightarrow$ <b><math>1.01 \times 10^6</math></b>
	Llama-3.3-70B-Instruct	7	0.63	5.41 $\rightarrow$ <b><math>3.86 \times 10^6</math></b>	8.65 $\rightarrow$ <b><math>2.18 \times 10^6</math></b>
Gemma	Gemma-2B	9	7.79	12.58 $\rightarrow$ <b><math>2.36 \times 10^{16}</math></b>	13.59 $\rightarrow$ <b><math>6.66 \times 10^{15}</math></b>
	Gemma-7B	3	1.01	9.98 $\rightarrow$ <b><math>6.25 \times 10^{21}</math></b>	11.41 $\rightarrow$ <b><math>1.32 \times 10^{20}</math></b>
DeepSeek-R1	Deepseek-R1-Distill-Qwen-1.5B	23	21.12	61.05 $\rightarrow$ <b><math>1.28 \times 10^3</math></b>	49.09 $\rightarrow$ <b><math>1.19 \times 10^3</math></b>
	DeepSeek-R1-Distill-Qwen-7B	3	1.28	34.57 $\rightarrow$ <b><math>5.59 \times 10^3</math></b>	35.62 $\rightarrow$ <b><math>3.70 \times 10^3</math></b>
	DeepSeek-R1-Distill-Llama-8B	3	1.33	17.56 $\rightarrow$ <b><math>2.81 \times 10^5</math></b>	23.43 $\rightarrow$ <b><math>2.55 \times 10^5</math></b>
	DeepSeek-R1-Distill-Qwen-14B	4	1.12	12.85 $\rightarrow$ <b><math>6.85 \times 10^3</math></b>	19.26 $\rightarrow$ <b><math>5.00 \times 10^3</math></b>
	DeepSeek-R1-Distill-Qwen-32B	28	3.58	9.36 $\rightarrow$ <b><math>8.07 \times 10^2</math></b>	15.14 $\rightarrow$ <b><math>7.40 \times 10^2</math></b>
	DeepSeek-R1-Distill-Llama-70B	3	0.27	7.86 $\rightarrow$ <b><math>2.21 \times 10^4</math></b>	11.34 $\rightarrow$ <b><math>1.84 \times 10^4</math></b>
Phi-3	Phi-3-mini-4k-Instruct	13	6.83	8.24 $\rightarrow$ <b><math>9.51 \times 10^4</math></b>	10.57 $\rightarrow$ <b><math>4.69 \times 10^4</math></b>
	Phi-3.5-mini-Instruct	13	6.83	8.46 $\rightarrow$ <b><math>2.67 \times 10^6</math></b>	10.71 $\rightarrow$ <b><math>1.02 \times 10^6</math></b>
Qwen2.5	Qwen2.5-0.5B-Instruct	3	5.90	18.18 $\rightarrow$ <b><math>1.76 \times 10^6</math></b>	22.34 $\rightarrow$ <b><math>1.08 \times 10^6</math></b>
	Qwen2.5-1.5B-Instruct	11	10.19	12.37 $\rightarrow$ <b><math>4.01 \times 10^2</math></b>	15.95 $\rightarrow$ <b><math>3.28 \times 10^2</math></b>
	Qwen2.5-3B-Instruct	5	2.89	11.24 $\rightarrow$ <b><math>3.34 \times 10^3</math></b>	14.28 $\rightarrow$ <b><math>2.52 \times 10^3</math></b>
	Qwen2.5-7B-Instruct	10	4.30	9.75 $\rightarrow$ <b><math>9.64 \times 10^4</math></b>	13.01 $\rightarrow$ <b><math>7.61 \times 10^4</math></b>
	Qwen2.5-14B-Instruct	4	1.13	7.69 $\rightarrow$ <b><math>2.39 \times 10^3</math></b>	11.12 $\rightarrow$ <b><math>1.59 \times 10^3</math></b>
	Qwen2.5-32B-Instruct	45	5.80	7.40 $\rightarrow$ <b><math>1.01 \times 10^2</math></b>	10.98 $\rightarrow$ <b><math>0.94 \times 10^2</math></b>
	Qwen2.5-72B-Instruct	3	0.26	10.23 $\rightarrow$ <b><math>2.24 \times 10^4</math></b>	10.72 $\rightarrow$ <b><math>1.40 \times 10^4</math></b>

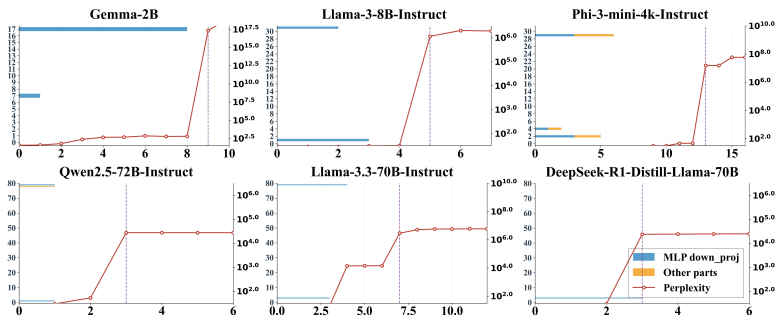
## Finding 2: Architectural Concentration

### Where are the critical neurons located?

- **Outer Layers:** Critical neurons are not uniformly distributed but exhibit a strong tendency to cluster in the outer layers.
- **Component Specificity:** They are almost exclusively located within the **MLP** `down_proj` components.
- **Interpretation:** This aligns with the "Super Weight" phenomenon. Outer `down_proj` layers act as extreme information bottlenecks, compressing high-dimensional representations back into the embedding space before final output generation.

# Finding 3: Phase Transition Behavior

- **Threshold Effect:** Masking partial subsets yields minimal impact; masking the *complete* set triggers a sudden, catastrophic perplexity surge.
- **Collective Circuitry:** These neurons operate as a highly interdependent computational unit rather than independent switches.



## Universal Capability Zeroing

We evaluated 70B-class models across 7 diverse benchmarks (MMLU-Pro, HumanEval, MATH, GPQA Diamond, IFEval, MGSM, SimpleQA).

- **Result:** After masking the identified critical neurons, accuracy/success rates dropped to **exactly 0.0000** across **all** tasks.
- This confirms that these sparse neurons control fundamental language processing capabilities, not just isolated or task-specific knowledge.

Model	Condition	MMLU-Pro	IFEval	GPQA-Diamond	HumanEval	MATH	MGSM	SimpleQA
Llama-3.3-70B-Instruct	Original	0.4441	0.4658	0.2879	0.2988	0.4420	0.9590	0.3704
	Masked	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
DeepSeek-R1-Distill-Llama-70B	Original	0.1631	0.2458	0.1869	0.1951	0.1420	0.9809	0.3104
	Masked	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
Qwen2.5-72B-Instruct	Original	0.2442	0.5268	0.3687	0.2866	0.3140	0.9080	0.2443
	Masked	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000

# Conclusion & Implications

- **Core Takeaway:** LLM capabilities are not as robustly distributed as assumed; they rely on an "Achilles' Heel" of ultra-sparse critical neurons.
- **Security Risks:** This extreme sparsity presents significant vulnerabilities, as minute, targeted interventions can trigger complete system failure.
- **Future Architecture:** Highlights the urgent need to rethink Transformer designs to enforce more uniform, redundant, and robust computational distribution across the network.

*Code available at: [github.com/qqqqqqqzx/The-Achilles-Heel-of-LLMs](https://github.com/qqqqqqqzx/The-Achilles-Heel-of-LLMs)*