

Spectral Scaling Laws in Language Models

How Effectively Do Feed-Forward Networks Use Their Latent Space?

Nandan Kumar Jha & Brandon Reagen (New York University)



AIW@ICML'25

Motivation: Why FFN Width Matters for LLMs?

- 1. Capacity and Computational Efficiency:
- Feedforward Networks (FFNs) contain **60-70%** of total parameters
- FFN FLOPs dominate in **<4K** context length regime
- FFNs store factual information and directly affect model's capacity
- 2. Architectural Diversity:
- Different models use different FFN width (GPT-2 & Pythia: **4**×, LLaMA: **2.67**×)
- 3. Interpretability Gap:
- Existing (Loss parameter) scaling laws ignore how FFN width is utilized

We need principled tools to analyze FFN capacity allocation and scaling efficiency

Our Approach: Spectral Rank Measures

Hard Spectral Rank (Participation Ratio)



Soft Spectral Rank (Shannon entropy rank)



How many dimensions are active in the eigenspectrum, relative to total variance

How uniformly variance is distributed in eigenspectrum (insensitive to spikes)

Key Findings: Asymmetric Spectral Scaling Laws LLaMA-70M (PreLN) LLaMA-130M (PreLN) LLaMA-250M (PreLN) LLaMA-250M(PreLN) - Dimension Scaling Laws LLaMA-130M(PreLN) - Dimension Scaling Laws LLaMA-70M(PreLN) - Dimension Scaling Laws --- HRank \propto D^{0.407 \pm 0.671} (R^2 = 0.268) --- HRank $\propto D^{0.604 \pm 0.411}$ (R²=0.684) --- HRank $\propto D^{0.451\pm0.778}$ (R²=0.251) --- SRank $\propto D^{1.069\pm0.292}$ (R²=0.930) --- SRank $\propto D^{0.872\pm0.353}$ (R²=0.859) 10^3 - --- SRank $\propto D^{0.879\pm0.490}$ (R²=0.763) 10^{3} 10^{3} sure



Asymmetric scaling pattern: Soft spectral rank exhibits better power law trend ($\beta \rightarrow 1$, $R^2 \rightarrow 1$) than hard spectral rank ($\beta \rightarrow 0.5$, $R^2 \rightarrow 0.5$)

Takeaways: Widening the FFN keeps adding low-energy directions (tail capacity) while the high-energy subspace reaches diminishing returns



500	1000	2000	3000	1000	2000	3000	5000	7000	1000	2000	3000	5000	7000
FFN Hidden Dimension (D)			FFN Hidden Dimension (D)				FFN Hidden Dimension (D)						

Takeaways: Hard spectral rank does not exhibit a sharp knee; a single sub-linear power law ($\beta \approx 0.5$) explains LLaMA-70M $\rightarrow 250M$

Impact of LayerNorm Positioning on Spectral Scaling Laws

	Р	reLN	Pe	ostLN	MixLN		
Model	Hard Rank	Soft Rank	Hard Rank	Soft Rank	Hard Rank	Soft Rank	
LLaMA-70M	0.451 ± 0.778	0.879 ± 0.490	0.556 ± 0.358	0.712 ± 0.273	0.593 ± 0.668	0.972 ± 0.477	
	$(R^2 = 0.251)$	$(R^2 = 0.763)$	$(R^2 = 0.706)$	$(R^2 = 0.872)$	$(R^2 = 0.440)$	$(R^2 = 0.805)$	
LLaMA-130M	0.604 ± 0.411	1.069 ± 0.292	0.521 ± 0.294	0.818 ± 0.372	0.626 ± 0.484	1.096 ± 0.484	
	$(R^2 = 0.684)$	$(R^2 = 0.930)$	$(R^2 = 0.758)$	$(R^2 = 0.829)$	$(R^2 = 0.626)$	$(R^2 = 0.837)$	

Takeaways: PostLN suppresses tail scaling; whereas MixLN achieves optimal capacity allocation by raising dominant-mode capacity and maintaining near-linear tail growth, delaying capacity saturation

Actionable Takeaways

Our spectral analysis demonstrate how widening the FFN primarily expands low-energy directions while dominant modes saturate earlier. We need architectural techniques (e.g., MixLN) to improve dominant mode energy without suppressing tail capacity

Contact: nj2049@nyu.edu