



Entropy-Guided Attention for Private LLMs

Nandan Kumar Jha & Brandon Reagen
New York University

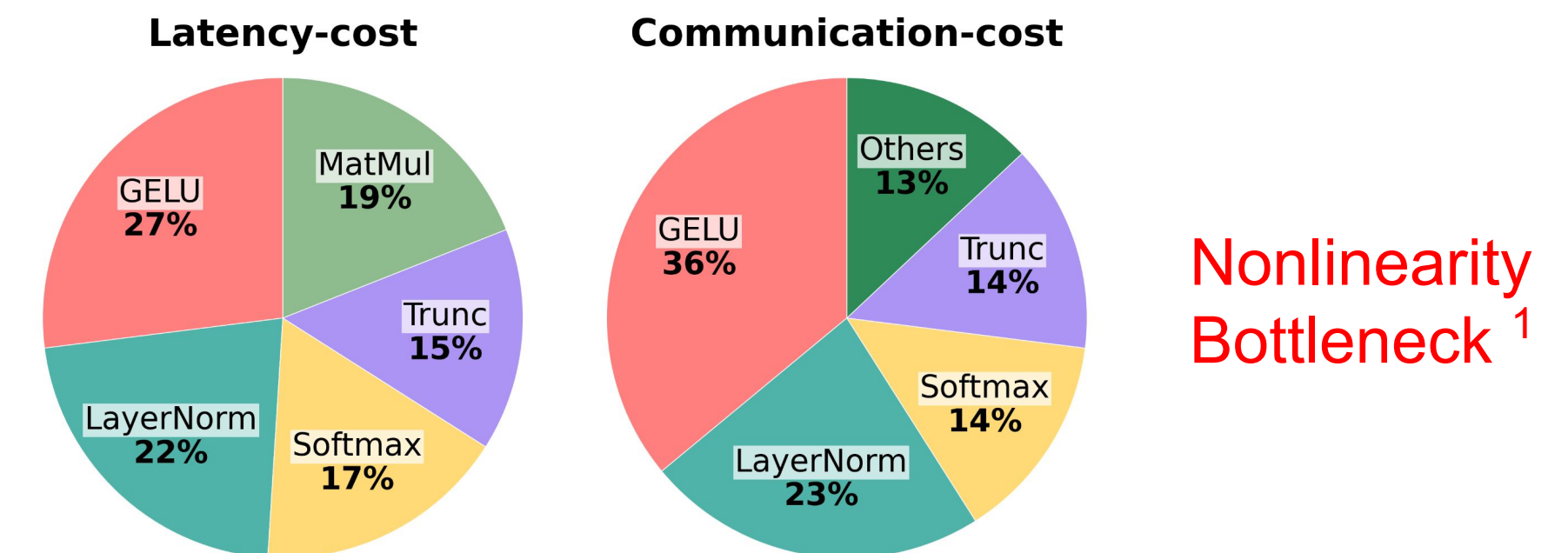


The 39th
Annual AAAI
Conference On
Artificial
Intelligence

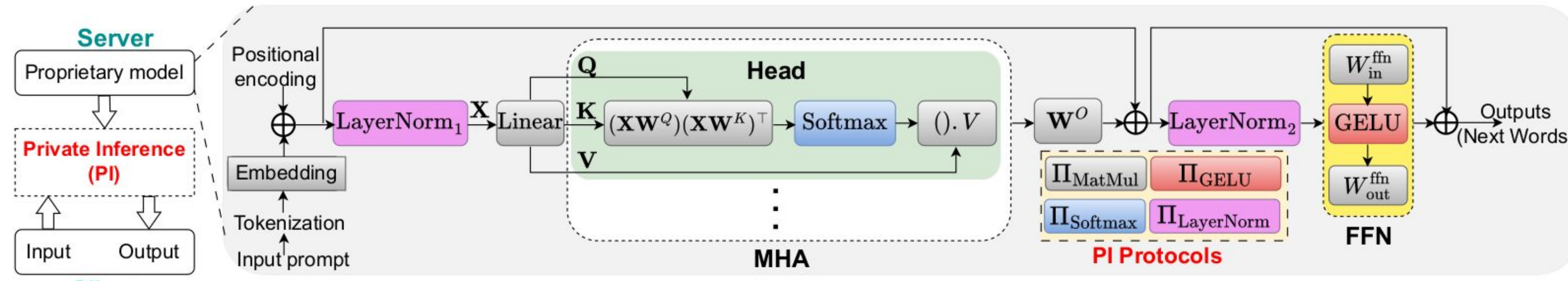
Private Inference on Large Language Models: An Introduction

Motivation and Challenges

High Latency & Communication Overheads: **8.2 minutes** and **25.3 GB** to generate a single token on GPT-2 (125M, T=128)

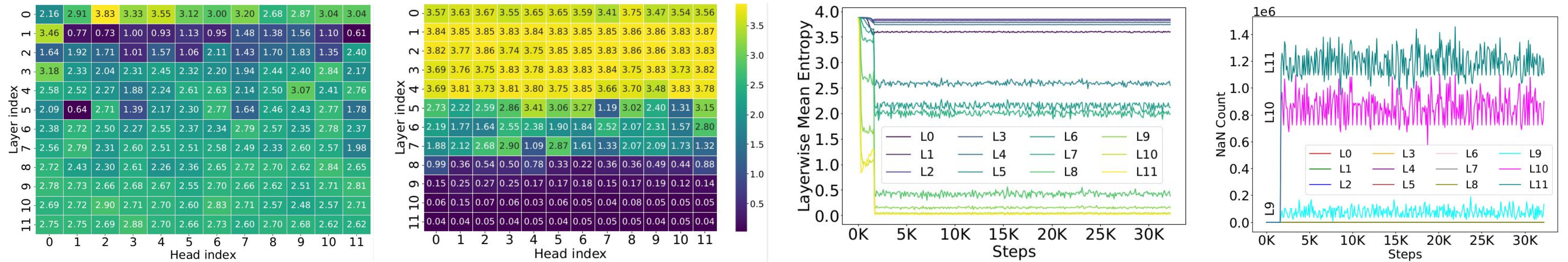


1. Hou et al., CipherGPT: Secure Two-Party GPT Inference



Computations performed directly on encrypted data, without seeing its content
Privacy of users' sensitive data is preserved while the server's model remains protected

Key Findings: The Absence of Nonlinearities in LLMs Leads to Entropic Overload & Entropy Collapse



Inference-Efficient Solutions to Prevent Entropy Collapse

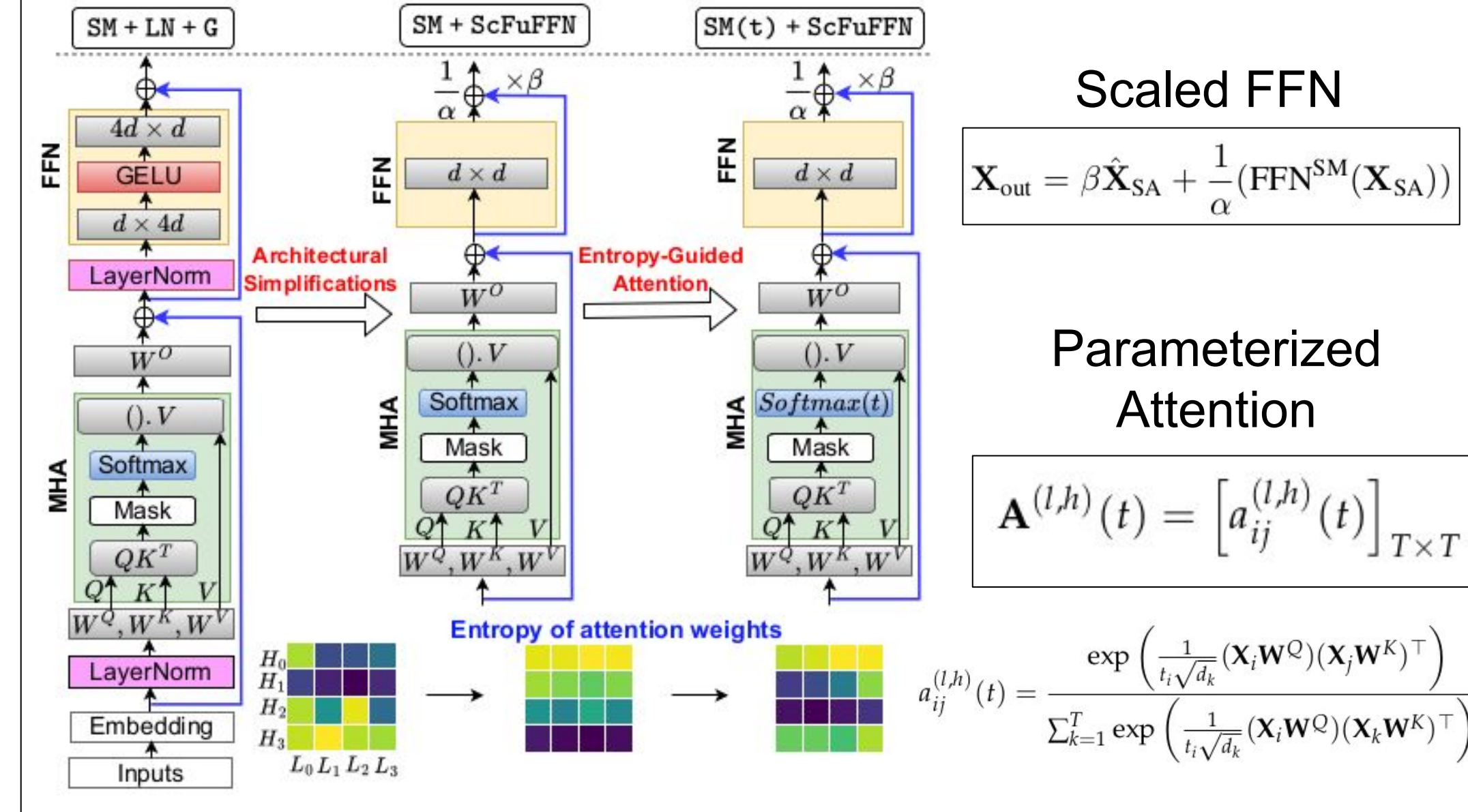
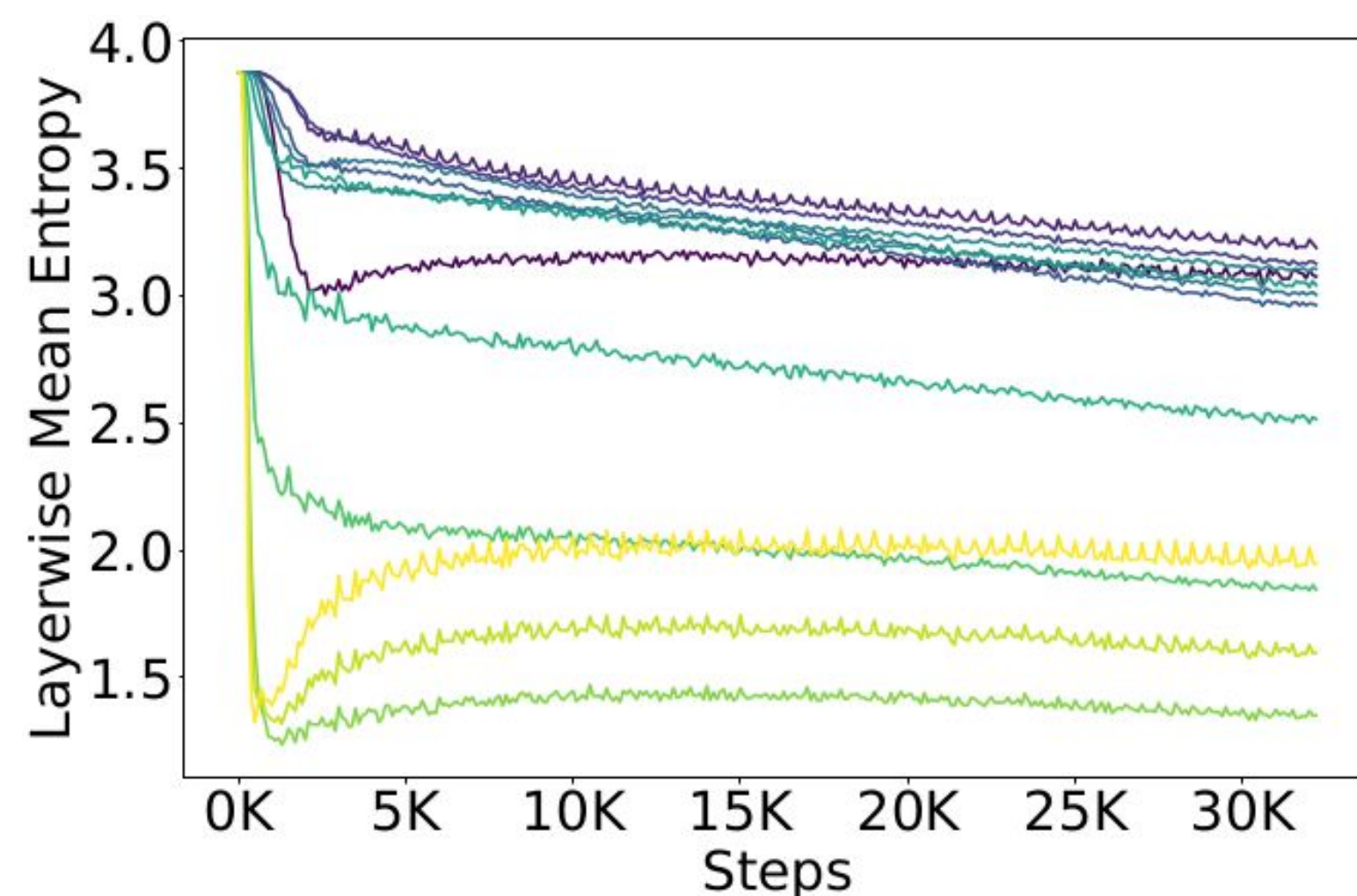
Entropy-Guided Private LLM Design

Solution 1: Weight Normalization in FFN

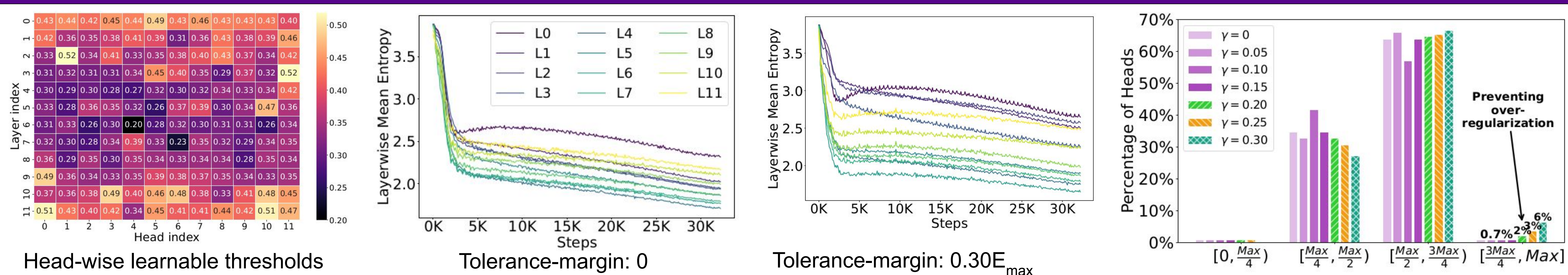
$$\text{FFN}_{\text{WNorm}}^{\text{SM}}(\mathbf{X}) = \left(\mathbf{X} \frac{\mathbf{V}_{\text{in}}^{\text{ffn}}}{\|\mathbf{V}_{\text{in}}\|_2} g_{\text{in}} \right) \frac{\mathbf{V}_{\text{out}}^{\text{ffn}}}{\|\mathbf{V}_{\text{out}}\|_2} g_{\text{out}}$$

Solution 2: Spectral Normalization in FFN

$$\text{FFN}_{\text{SNorm}}^{\text{SM}}(\mathbf{X}) = \left(\mathbf{X} \frac{\mathbf{W}_{\text{in}}^{\text{ffn}}}{\sigma(\mathbf{W}_{\text{in}}^{\text{ffn}})} \right) \frac{\mathbf{W}_{\text{out}}^{\text{ffn}}}{\sigma(\mathbf{W}_{\text{out}}^{\text{ffn}})}$$



Key Innovation in Entropy Regularization Scheme: Learnable Threshold and Tolerance Margin



Experimental Results: GPT-2 (L=12, H=12, d=768)

Network Arch.	PPL	#Nonlinear Ops	#FLOPs				Savings	
			FFN	Attn.	Comm. (GB)	Lat. (min.)	Comm.	Lat.
Baseline SM + LN + G	2.69	SM:144 × ℝ ^{128×128} LN:24 × ℝ ^{128×768} G:12 × ℝ ^{128×3072}	14.5B	7.7B	25.32	8.21	1×	1×
			SM:144 × ℝ ^{128×128} LN:24 × ℝ ^{128×768} R:12 × ℝ ^{128×3072}	14.5B	7.7B	9.44	6.06	2.68×
SM + ScFuFFN	3.48	SM:144 × ℝ ^{128×128}	1.8B	7.7B	6.43	4.76	3.94×	1.72×
EReg(SM(t) + ScFuFFN)	3.21	SM:144 × ℝ ^{128×128}	1.8B	7.7B	6.43	4.76	3.94×	1.72×

CodeParrot Dataset (2.1B Tokens, T=128)

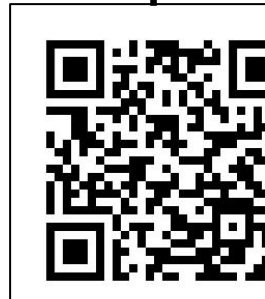
Network Arch.	Eval PPL			#Nonlinear Ops	#FLOPs			
	1.2B	2.4B	4.8B		FFN	Attn.	Comm. (GB)	Lat. (min.)
Baseline SM + LN + G	25.71	23.32	21.29	SM:144 × ℝ ^{512×512} LN:24 × ℝ ^{512×768} G:12 × ℝ ^{512×3072}	58.0B	36.2B	145.24	30.74
	SM + LN + R	26.06	23.55	21.58	SM:144 × ℝ ^{512×512} LN:24 × ℝ ^{512×768} R:12 × ℝ ^{512×3072}	58.0B	36.2B	81.71
SM + ScFuFFN	33.77	30.82	28.59	SM:144 × ℝ ^{512×512}	7.3B	36.2B	69.68	19.44
EReg(SM(t) + ScFuFFN)	31.54	28.70	26.55	SM:144 × ℝ ^{512×512}	7.3B	36.2B	69.68	19.44

Languini Book Dataset (1.2B to 4.8B Tokens, T=512)

Conclusion and Key Takeaways

While nonlinearities are essential in LLMs, **strategically** applied entropy regularization and FFN normalization can prevent entropy collapse and entropic overload in (Softmax-only) private LLMs.

Paper



Code



Contact:
nj2049@nyu.edu