

# cMoLLM at Scale:

## Horizontal Scaling Laws for Convolutionally-Gated Mixture-of-LLMs

Xin Yang · Yemin Wang · Mingda Liu · Letian Li · Shuaishuai Cao · Zhengxiao He · Ryan Dong

Seoul, South Korea · 2026

Dynamic convolution view of MoE

Pipeline-level scaling

Stable soft routing

**01**

## 1. Motivation

Dense scaling and limits of prior MoE-style methods

---

**02**

## 2. Key Insight

MoE layers as dynamic  $1 \times 1$  convolutions

---

**03**

## 3. cMoLLM Method

Soft kernel mixing over end-to-end streams

---

**04**

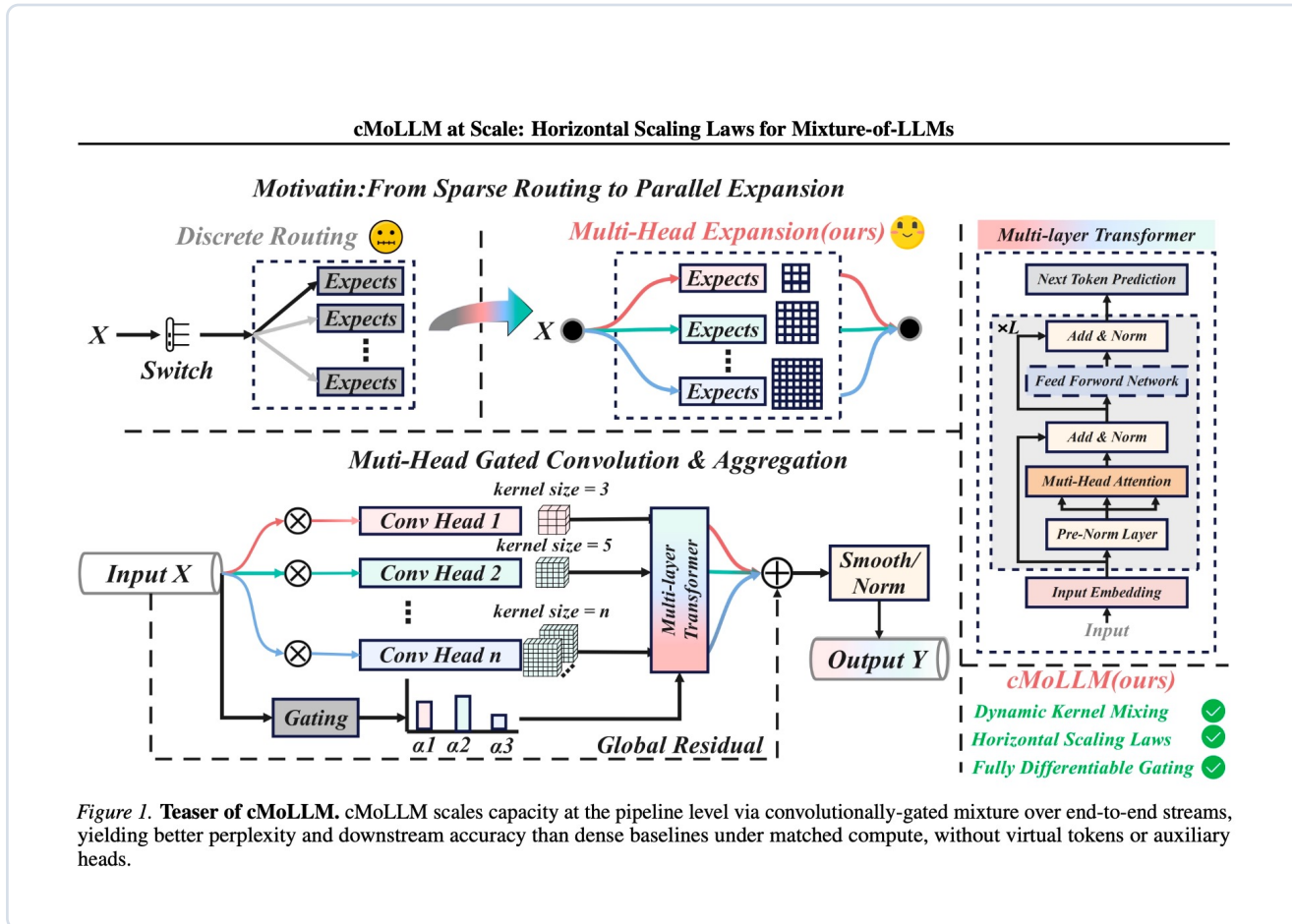
## 4. Results & Scaling

Ablations, scaling laws, and future directions

---

# Motivation: Capacity Is Still Tied to Compute

Dense Transformers activate every parameter for every token, making scaling increasingly expensive.



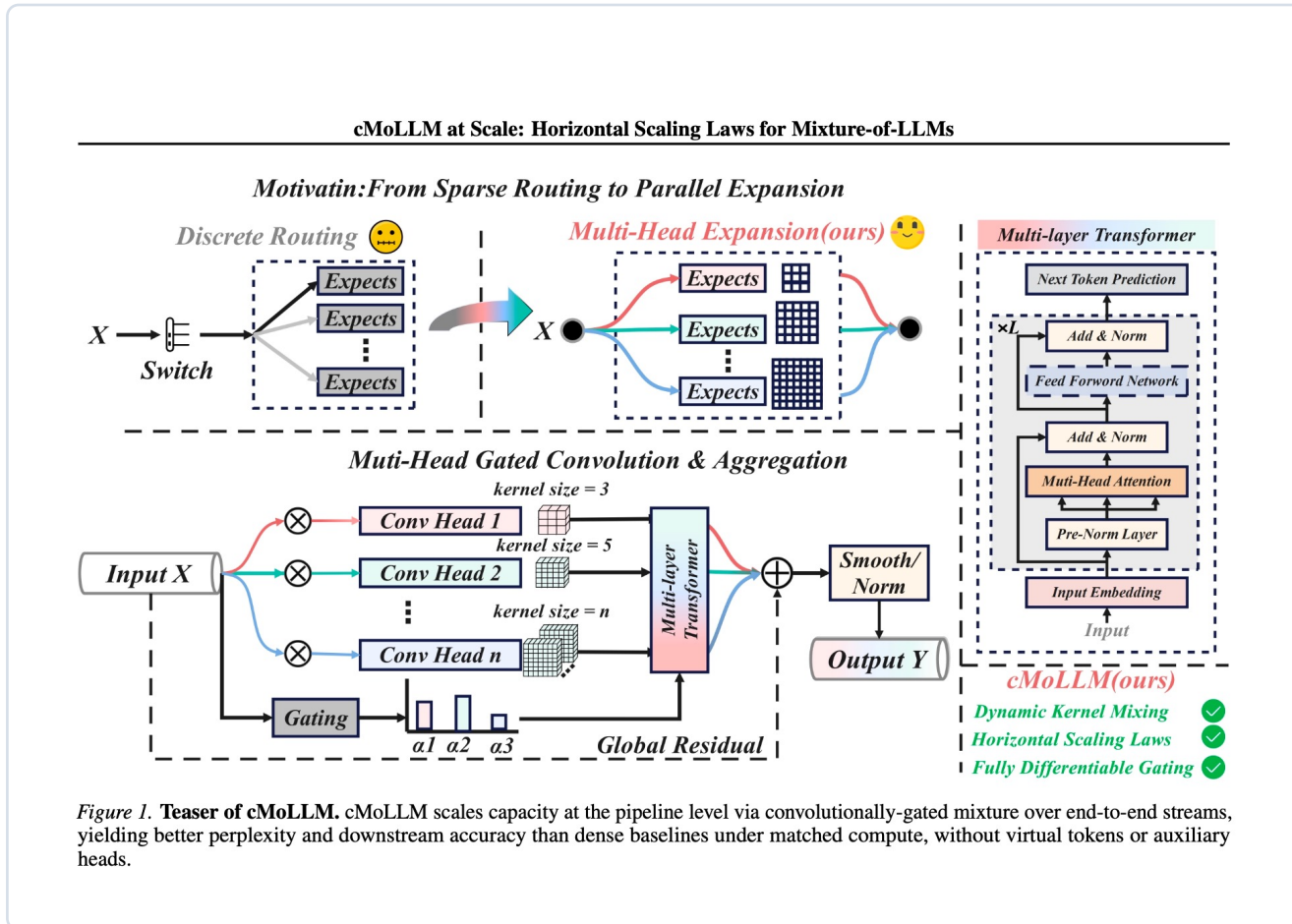
## Core Problem

- Dense LLMs couple capacity and activated computation.
- Classical MoE usually improves capacity only inside FFN blocks.
- Sparse Top-K routing can introduce expert collapse and brittle utilization.
- The paper asks: can we scale capacity horizontally without unstable routing?

Desired direction: pipeline-level mixture with stable, hardware-friendly routing.

# Why Prior Pipeline-Level Scaling Is Not Enough

ParaScale and AltUp are useful baselines, but each brings extra structural assumptions or optimization risks.



Observed Gap

## The missing design point

- No virtual tokens: avoid sequence-length overhead.
- No auxiliary branch: keep the pipeline simple and adaptive.
- No sparse Top-K: reduce routing discontinuity and collapse risk.
- Keep computation close to dense Transformer cost for small N.

- cMoLLM targets this gap with dynamic convolutional routing.

# Key Insight: MoE as Dynamic 1×1 Convolution

Experts become kernels; routing becomes input-conditioned kernel aggregation.

## Theorem 4.1

For linear experts:

$$\text{MoE}(\mathbf{x}) = \sum_k \mathbf{g}_k(\mathbf{x}) \mathbf{E}_k(\mathbf{x})$$

If  $\mathbf{E}_k(\mathbf{x}) = \mathbf{x} \mathbf{W}_k^\top$ , then

$$\text{MoE}(\mathbf{x}) = \text{Conv}_{1 \times 1}(\mathbf{x}; \tilde{\mathbf{K}}(\mathbf{x}))$$

$$\tilde{\mathbf{K}}(\mathbf{x}) = \sum_k \mathbf{g}_k(\mathbf{x}) \mathbf{W}_k$$

Interpretation: dynamic convolution provides a continuous, differentiable lens for conditional computation.

$$\text{MoE}(\mathbf{x}) = \sum_{k=1}^N g_k(\mathbf{x}) E_k(\mathbf{x}), \quad (2)$$

where  $g(\mathbf{x}) \in \mathbb{R}^N$  denotes the routing weights produced by the gate. In practice, only the top- $K$  entries of  $g(\mathbf{x})$  are nonzero, yielding sparse computation. Each expert  $E_k$  is typically a two-layer FFN with its own parameters.

### 3.4. Pointwise Convolution as a Linear Transform

A pointwise ( $1 \times 1$ ) convolution applies the same linear map across token positions. Given an input sequence  $\mathbf{X} \in \mathbb{R}^{L \times d_{in}}$  and a kernel  $\mathbf{K} \in \mathbb{R}^{d_{out} \times d_{in}}$ , the operation is given by

$$\text{Conv}_{1 \times 1}(\mathbf{X}; \mathbf{K}) = \mathbf{X} \mathbf{K}^\top, \quad (3)$$

which corresponds to applying an identical token-wise linear projection from  $\mathbb{R}^{d_{in}}$  to  $\mathbb{R}^{d_{out}}$ . This observation serves as a key building block for our subsequent analysis.

## 4. Method: cMoLLM

We now present our main theoretical result and the cMoLLM architecture.

### 4.1. MoE as Dynamic Convolution: A Formal Equivalence

$$\begin{aligned} \sum_{k=1}^N g_k(\mathbf{x}) \cdot (\mathbf{x} \mathbf{W}_k^\top) &= \mathbf{x} \left( \sum_{k=1}^N g_k(\mathbf{x}) \mathbf{W}_k \right) \\ &= \text{Conv}_{1 \times 1} \left( \mathbf{x}; \sum_{k=1}^N g_k(\mathbf{x}) \mathbf{W}_k \right). \end{aligned} \quad (6)$$

The weighted sum of expert weight matrices forms a dynamic kernel that varies with  $\mathbf{x}$ . (Classical MoE uses Top- $K$  routing so only  $K$  terms are nonzero; we use soft mixing over all streams.) See Section A for extension to nonlinear experts.

**Corollary 4.2** (Nonlinear Two-Layer Experts). *Under Theorem 3.2, consider two-layer experts of the form  $E_k(\mathbf{x}) = \sigma(\mathbf{x} \mathbf{W}_{k,1}^\top) \mathbf{W}_{k,2}^\top$  with routing weights  $\{g_k(\mathbf{x})\}_{k=1}^N$  satisfying  $\sum_k g_k(\mathbf{x}) = 1$ . Then the MoE layer can still be written as a dynamic  $1 \times 1$  convolution with an input-dependent effective kernel that incorporates a data-dependent mask induced by  $\sigma$ ; see Section A for details.*

Assumptions 3.1 and 3.2 are not needed for Theorem 4.1 itself, but they will be used in our later complexity, stability, and toy-model analyses that build on this equivalence.

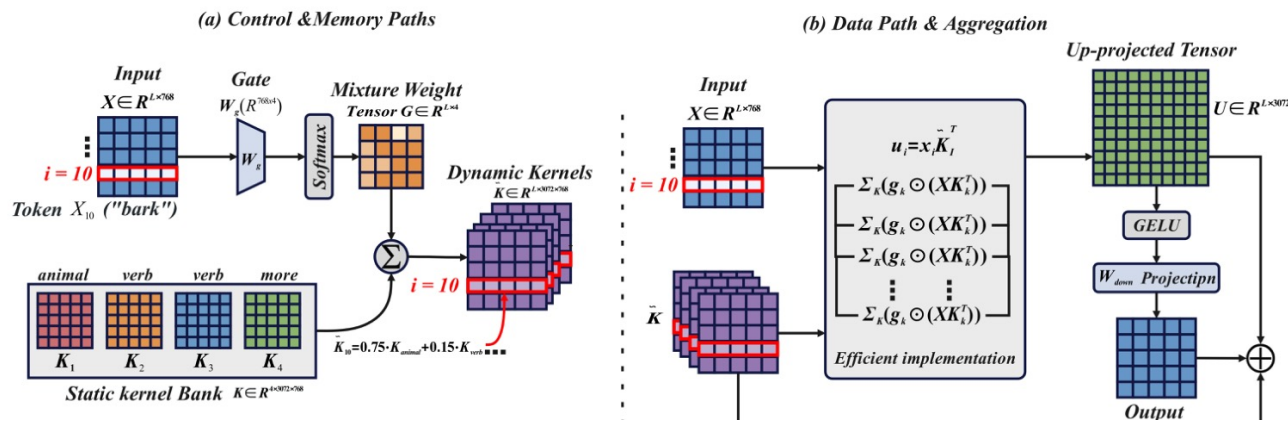
**Remark 4.3** (Key Differences from Standard Convolution). While Theorem 4.1 reveals structural similarity, standard convolutional layers differ from explicit MoE in three aspects:

1. **Static vs. Learnable Gates.** In CNN, activation func

# cMoLLM Architecture

A convolutionally-gated mixture over end-to-end LLM streams.

## cMoLLM at Scale: Horizontal Scaling Laws for Mixture-of-LLMs



### Design Recipe

- Keep the Transformer backbone.
- Replace FFN stack with stream-wise kernels.
- Use a lightweight gate to produce soft mixture weights.
- Apply dynamic  $1 \times 1$  convolution through grouped / pointwise primitives.

No virtual tokens

No auxiliary heads

No Top-K

No low-rank

# Forward Pass and Training Strategy

A single cMoLLM block uses soft weights to mix per-stream kernels before GELU and down-projection.

---

## Algorithm 1 Forward pass of a Transformer layer with cMoLLM

---

- 1: **Input:**  $\mathbf{X} \in \mathbb{R}^{L \times d}$ , number of streams  $N$ , kernels  $\{\mathbf{K}_k\}_{k=1}^N$ ,  $\mathbf{W}_{\text{down}}$
  - 2: **Self-Attention:**  $\mathbf{H} \leftarrow \text{SelfAttention}(\mathbf{X})$
  - 3: **Gating logits:**  $\mathbf{G} \leftarrow \mathbf{H}\mathbf{W}_g + \mathbf{b}_g \in \mathbb{R}^{L \times N}$
  - 4: **Soft mixture weights:**  $g_{i,k} \leftarrow \text{softmax}(\mathbf{G}_{i,:})_k$  for each token  $i$ , stream  $k$
  - 5: **Mixed kernel:**  $\tilde{\mathbf{K}}_i \leftarrow \sum_{k=1}^N g_{i,k} \mathbf{K}_k$
  - 6: **Dynamic 1×1 conv:**  $\mathbf{U}_i \leftarrow \mathbf{H}_i \tilde{\mathbf{K}}_i^\top$
  - 7: **Nonlinearity & down-proj:**  $\mathbf{Y}_i \leftarrow \mathbf{W}_{\text{down}} \sigma(\mathbf{U}_i)$
  - 8: **Output:**  $\mathbf{Y}$  (residual + norm as usual)
- 

### Forward Pass

- Self-attention produces hidden states  $\mathbf{H}$ .
- A softmax gate estimates token-wise stream weights  $g_{i,k}$ .
- The mixed kernel  $\tilde{\mathbf{K}}_i = \sum_k g_{i,k} \mathbf{K}_k$  is applied as dynamic 1×1 convolution.

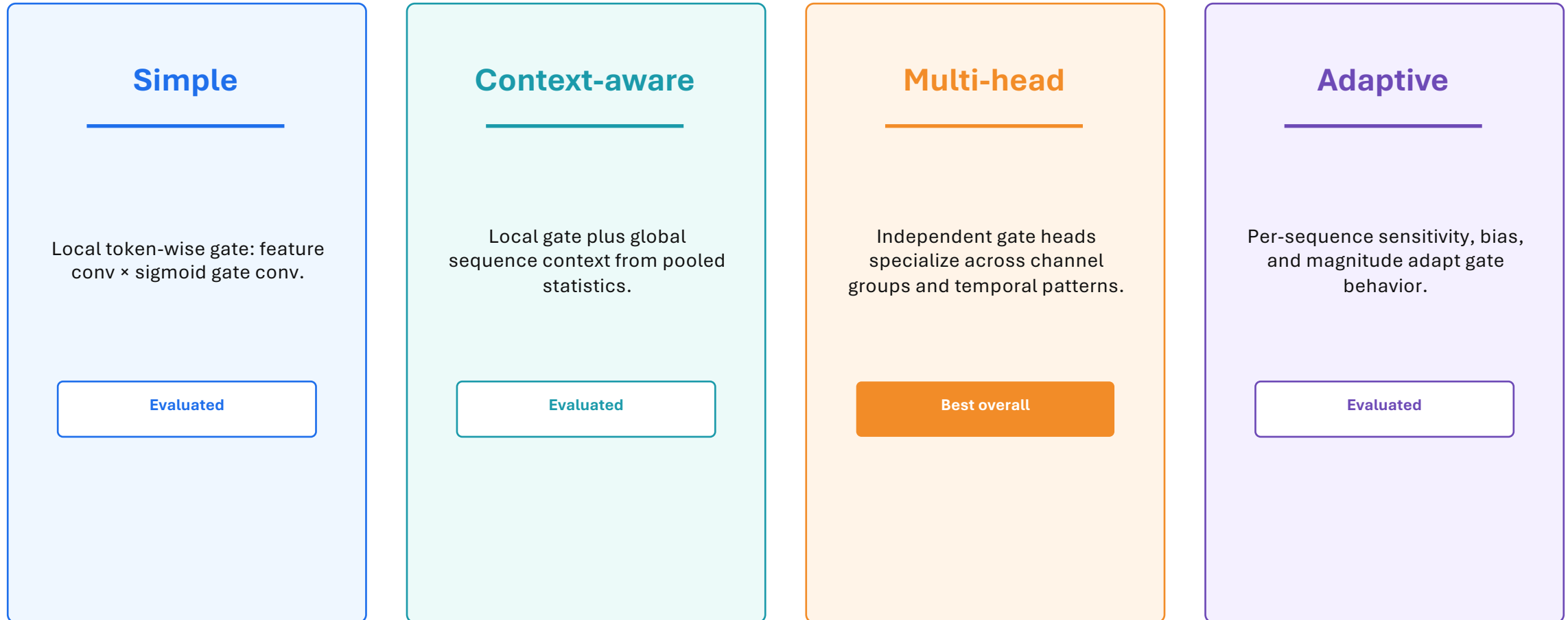
### Load Balancing

$$\mathbf{L} = \text{LLM} + \text{Lbal} \quad \text{Lbal} = \alpha \cdot \mathbf{N} \cdot \sum_k p_k^2$$

- Discourages collapse onto a few streams.
- Keeps routing fully differentiable.
- No kernel regularization, progressive sparsification, or low-rank terms.

# Gating Variants

The paper evaluates four convolutional gates; multi-head gives the strongest overall trade-off.



**Result pattern: increasing stream count generally improves perplexity; multi-head remains stable at N = 8.**

# Complexity: Horizontal Scaling Law

For bounded  $N$ , capacity grows through streams while per-token compute remains a constant-factor expansion.

## cMoLLM at Scale: Horizontal Sca

### 4.5.4. ADAPTIVE GATED CONVOLUTION

An adaptive network (global pooling plus a two-layer MLP) outputs three bounded parameters (sensitivity, bias, magnitude) from global input statistics. The gate is sigmoid(scaled gate-conv output + bias) times magnitude. Gating behavior is adjusted per sequence (e.g., dynamic range or complexity).

### 4.6. Training Strategies

**Load Balancing Loss.** To encourage balanced use of streams, we add an auxiliary loss (Fedus et al., 2022):

$$\mathcal{L}_{\text{bal}} = \alpha \cdot N \cdot \sum_{k=1}^N p_k^2, \quad (20)$$

where  $p_k = \frac{1}{L} \sum_{i=1}^L g_k(\mathbf{x}_i)$  is the average gating probability for stream  $k$  over the sequence. Minimizing  $\sum_k p_k^2$  discourages collapse onto a few streams. We use  $\alpha = 0.01$  in experiments.

**Total Loss.** The total training objective is  $\mathcal{L} = \mathcal{L}_{\text{LM}} + \mathcal{L}_{\text{bal}}$ ;

$$\text{FLOPs}_{\text{cMoLLM}} \approx (N + 1) d d_{\text{ff}} + \text{FLOPs}_{\text{gate}}$$

Horizontal scaling: effective capacity grows with the number of streams while per-token cost stays close to the dense baseline for small  $N$ .

# Experimental Setup

GPT-2-style language modeling and downstream evaluation under matched training settings.

<b>Base architecture</b>	GPT-2-style Transformer
<b>Training data</b>	FineWeb large-scale curated web corpus
<b>Model scales</b>	85M, 350M, 760M parameters
<b>Sequence length</b>	4096
<b>FFN dimension</b>	3072
<b>Learning rate</b>	$6 \times 10^{-5}$
<b>Metrics</b>	validation loss, PPL, GLUE, SQuAD v2
<b>Reporting</b>	3 random seeds; mean $\pm$ std

Ablation: gating variants  $\times$  stream counts  $N \in \{1, 2, 4, 8\}$

Scaling: 85M  $\rightarrow$  350M  $\rightarrow$  760M

# Main Ablation Results: 85M Scale

Multi-head gating with N = 8 achieves the best loss, PPL, and SQuAD v2 in Table 4.

**cMoLLM at Scale: Horizontal Scaling Laws for Mixture-of-LLMs**

Table 4. Full experimental results (3 seeds; mean  $\pm$  std): validation loss, perplexity (PPL), GLUE (%), SQuAD v2 (%). “conv” = gating variant;  $n$  = streams. **Bold** = best (SOTA) in column.

conv	$n$	Loss $\downarrow$	PPL $\downarrow$	GLUE (%) $\uparrow$	SQuAD v2 (%) $\uparrow$
base	–	3.23 $\pm$ 0.03	25.40 $\pm$ 0.65	42.45 $\pm$ 0.99	50.07 $\pm$ 0.65
simple	1	3.04 $\pm$ 0.03	20.94 $\pm$ 0.22	44.16 $\pm$ 0.56	51.13 $\pm$ 0.54
simple	2	3.01 $\pm$ 0.02	20.39 $\pm$ 0.56	44.62 $\pm$ 1.02	52.08 $\pm$ 0.95
simple	4	2.97 $\pm$ 0.04	19.53 $\pm$ 0.49	44.92 $\pm$ 0.87	55.71 $\pm$ 1.30
simple	8	3.00 $\pm$ 0.03	20.24 $\pm$ 0.12	44.68 $\pm$ 0.33	55.95 $\pm$ 0.51
context_aware	1	3.04 $\pm$ 0.03	20.94 $\pm$ 0.41	44.16 $\pm$ 0.42	51.12 $\pm$ 1.28
context_aware	2	2.98 $\pm$ 0.03	19.83 $\pm$ 0.55	46.79 $\pm$ 1.21	52.13 $\pm$ 1.13
context_aware	4	2.98 $\pm$ 0.02	19.78 $\pm$ 0.12	47.12 $\pm$ 0.25	55.64 $\pm$ 0.87
context_aware	8	2.96 $\pm$ 0.01	19.31 $\pm$ 0.68	<b>48.98<math>\pm</math>0.64</b>	56.11 $\pm$ 0.73
multi_head	1	3.04 $\pm$ 0.01	20.93 $\pm$ 0.22	44.16 $\pm$ 1.18	51.18 $\pm$ 0.66
multi_head	2	2.99 $\pm$ 0.05	19.90 $\pm$ 0.28	47.78 $\pm$ 0.45	52.98 $\pm$ 1.35
multi_head	4	2.96 $\pm$ 0.04	19.44 $\pm$ 0.41	48.02 $\pm$ 0.79	56.37 $\pm$ 1.56
multi_head	8	<b>2.93<math>\pm</math>0.02</b>	<b>18.82<math>\pm</math>0.37</b>	48.82 $\pm$ 1.01	<b>57.06<math>\pm</math>0.52</b>
adaptive	1	3.04 $\pm$ 0.05	20.94 $\pm$ 0.46	44.16 $\pm$ 0.98	51.28 $\pm$ 0.88
adaptive	2	3.00 $\pm$ 0.02	20.14 $\pm$ 0.26	46.64 $\pm$ 0.43	52.43 $\pm$ 1.45
adaptive	4	2.97 $\pm$ 0.03	19.68 $\pm$ 0.71	47.03 $\pm$ 0.36	55.89 $\pm$ 1.39
adaptive	8	2.97 $\pm$ 0.03	19.58 $\pm$ 0.65	47.81 $\pm$ 0.75	56.74 $\pm$ 0.72

Best Row

2.93

Loss  $\downarrow$

18.82

PPL  $\downarrow$

48.82

GLUE  $\uparrow$

57.06

SQuAD v2  $\uparrow$

**Best cMoLLM variant: multi-head, N = 8**

Key message: soft stream routing gives consistent gains over the dense base under the same evaluation protocol.

# Cross-Scale Results

The advantage of multi-head cMoLLM is preserved from 85M to 760M parameters.

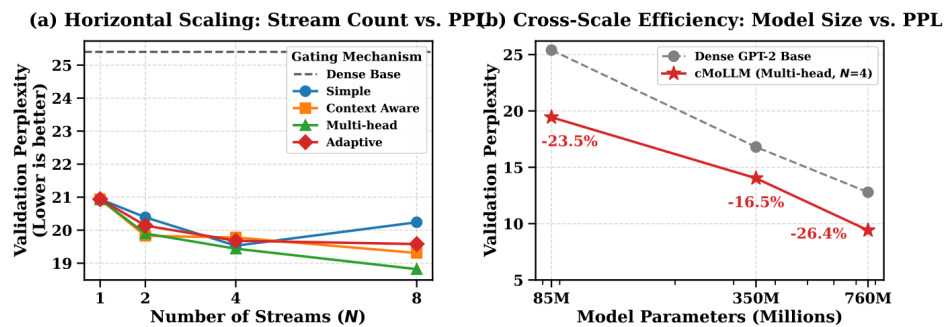
Table 5. Scaling experiments (3 seeds; mean  $\pm$  std) across model sizes (Table 2). mh = multi-head cMoLLM; base = dense baseline. **Bold** = SOTA in column.

type	Loss $\downarrow$	PPL $\downarrow$	GLUE $\uparrow$	SQuAD v2 $\uparrow$
base-85M	3.23 $\pm$ 0.03	25.40 $\pm$ 0.65	42.45 $\pm$ 0.99	50.07 $\pm$ 0.65
base-350M	2.82 $\pm$ 0.04	16.80 $\pm$ 0.47	50.68 $\pm$ 1.08	59.24 $\pm$ 1.24
base-760M	2.55 $\pm$ 0.05	12.79 $\pm$ 0.28	55.74 $\pm$ 0.46	62.91 $\pm$ 0.74
mh-85M	2.96 $\pm$ 0.03	19.44 $\pm$ 0.56	48.02 $\pm$ 0.84	56.37 $\pm$ 0.54
mh-350M	2.64 $\pm$ 0.03	14.03 $\pm$ 0.37	53.82 $\pm$ 1.03	61.08 $\pm$ 1.26
mh-760M	<b>2.24<math>\pm</math>0.02</b>	<b>9.41<math>\pm</math>0.45</b>	<b>58.45<math>\pm</math>0.77</b>	<b>64.56<math>\pm</math>1.32</b>

## Cross-scale takeaway

- At every reported scale, cMoLLM improves loss and perplexity over the dense baseline.
- The 760M cMoLLM model reaches PPL 9.41 versus 12.79 for the dense base.
- Downstream GLUE and SQuAD v2 accuracy also improve at all three model scales.

cMoLLM at Scale: Horizontal Scaling Laws for Mixture-of-LLMs

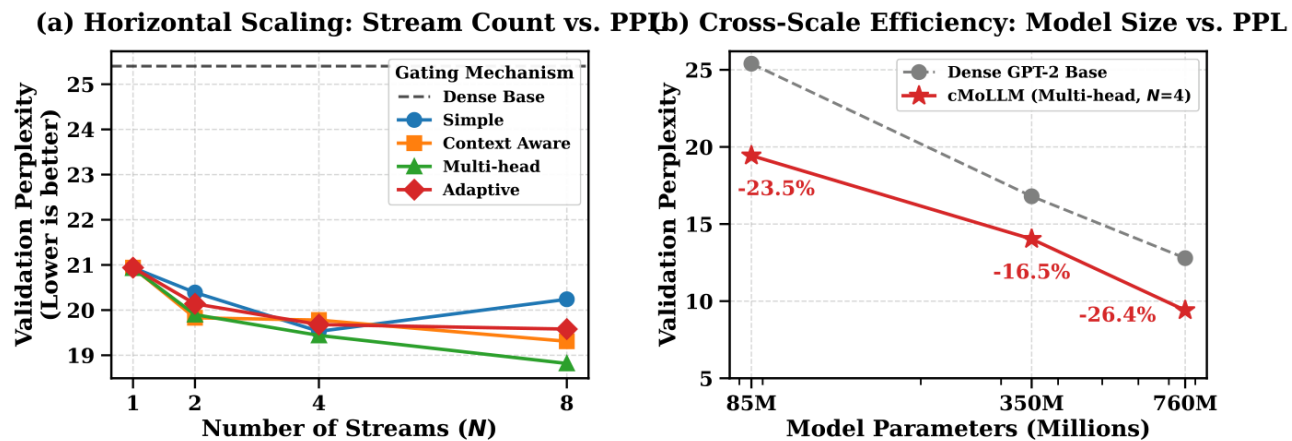


Scaling claim: cMoLLM scales favorably with both stream count and parameter count.

# Scaling Behavior: Stream Count and Model Size

Horizontal scaling improves perplexity as N increases, while cross-scale gains remain stable.

cMoLLM at Scale: Horizontal Scaling Laws for Mixture-of-LLMs



## Interpretation

- N = 1 → 4/8 reduces validation PPL across most gating variants.
- Multi-head remains stable at high stream count.
- The cMoLLM curve stays below the dense baseline at all tested model sizes.
- Empirical evidence matches the proposed horizontal scaling law.

# Summary: What cMoLLM Buys

A compact view of the paper's theoretical, architectural, and empirical contributions.

## Theoretical Lens

MoE-style routing is equivalent to input-conditioned dynamic  $1 \times 1$  convolution.

---

## Architectural Design

End-to-end streams are softly mixed via convolutional kernels rather than sparse Top-K selection.

---

## Empirical Behavior

Better PPL, GLUE, and SQuAD with stable utilization and favorable scaling.

---

**cMoLLM: more effective capacity under matched compute, with fully differentiable routing and no virtual tokens / auxiliary heads / Top-K.**

# Limitations, Future Work, and Q&A

The method is promising, but larger-scale validation remains the next key step.

## Limitations

- Current experiments are at GPT-2 scale up to approximately 760M parameters.
- 7B+ validation and distributed training need to be tested.
- The equivalence is explored mainly for mixture layers; attention-side extensions remain open.

## Future Work

- Scale to larger models and hardware clusters.
- Extend MoE-convolution equivalence to attention.
- Combine with retrieval or other conditional computation mechanisms.

# Thank you

# Questions?

cMoLLM at Scale  
Horizontal Scaling Laws for Convolutionally-Gated  
Mixture-of-LLMs

