

Uncertainty-Aware Routing for Principled Alignment with MoE Dynamics

Anonymous ACL submission

Abstract

Mixture-of-Experts (MoE) is a cornerstone for scaling LLMs, yet its training dynamics remain poorly understood, often leading to sub-optimal specialization. Moving beyond static routing, we present a systematic study of the MoE lifecycle using **Helmholtz Free Energy** and **Router Entropy**. We identify a universal **Three-Stage Phase Transition**—Exploration, Symmetry Breaking, and Stabilization—marked by an Energy “Climb” and Plateau. This reflects **Frustrated Exploration**, caused by structural interference between specialization drives and uniformity constraints. To address this, we propose **Uncertainty-Aware Routing (UAR)**, which aligns routing with the model’s epistemic state via: (1) **Evidence-Triggered Expansion**, increasing active experts for high-energy tokens, and (2) **Epistemic Masking**, applying load-balancing only in high-uncertainty regimes to shield mature experts. Experiments confirm UAR reduces perplexity and improves expert distinctiveness, offering a principled path toward thermodynamically aligned computation.

1 Introduction

The progress of Large Language Models (LLMs) has largely followed empirical scaling laws (Kaplan et al., 2020; Hoffmann et al., 2022; OpenAI, 2023), linking performance gains to increased parameter counts. As models approach the trillion-parameter regime, however, dense architectures impose unsustainable computational costs (Brown et al., 2020; Liu et al., 2023). Mixture-of-Experts architectures have therefore become the dominant strategy for efficient scaling (Shazeer et al., 2017; Jiang et al., 2024; Yang et al., 2025). By activating only a sparse subset of experts per token, MoE decouples model capacity from inference cost and enables training at unprecedented scale.

Despite their success, MoE models remain fragile to train. Their performance is highly sensitive

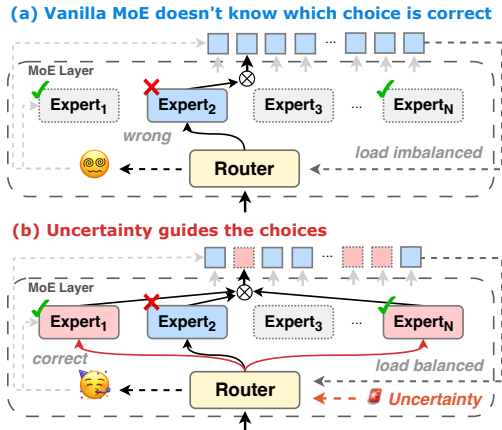


Figure 1: (a) The vanilla MoE router often selects incorrect experts and suffers from load imbalance due to a lack of guidance. (b) By incorporating uncertainty signals, the router is guided to select the correct experts, achieving higher accuracy and better load balancing.

to routing behavior, auxiliary loss coefficients, and training schedules (Lepikhin et al., 2020; Wei et al., 2024). Prior work has focused primarily on mitigating surface-level failures such as load imbalance, typically through auxiliary losses or heuristic routing constraints (Fedus et al., 2022; Du et al., 2022). These methods implicitly assume that routing is a static classification problem. What they overlook is that routing decisions evolve throughout training, and that this evolution fundamentally shapes expert learning (Zhou et al., 2022; Huang et al., 2024; Feng et al., 2025). Without a principled understanding of router dynamics, the learning processes of the router and experts remain poorly synchronized.

To bridge this gap, we conduct a systematic empirical study of MoE dynamics through a thermodynamic lens, analyzing the interaction between router entropy (aleatoric uncertainty) (Shannon, 1948; Bishop and Nasrabadi, 2006; Lepikhin et al., 2020) and Helmholtz Free Energy (epistemic uncertainty) (Helmholtz, 1982). Crucially, our anal-

ysis identifies a counter-intuitive **Energy Climb** followed by a high-level plateau. Rather than decreasing monotonically, rises sharply during the transitional phase, signaling a state of **Frustrated Exploration**. This phenomenon stems from a conflict between the main objective’s drive for specialization and the auxiliary loss’s push for uniformity, which traps the router in extreme epistemic uncertainty and prevents belief consolidation.

This thermodynamic friction manifests through several consistent mechanisms across model scales. We observe that strong auxiliary losses act as a thermal pressure suppressing logit magnitudes, while a distinct symmetry-breaking transition (5k–15k steps) marks the differentiation of token difficulty. Deeper layers, however, maintain a **Cognitive Vacuum**—a persistent energy gap where representation complexity outpaces evidence accumulation. Ultimately, these dynamics expose a fundamental **Diversity–Specialization Dilemma**: aggressive load balancing enforces uniformity at the cost of distinctiveness, collapsing specialists into redundant generalists and obstructing expert maturation under high-energy conditions.

Guided by these insights, we propose **Uncertainty-Aware Routing (UAR)**, a framework that explicitly aligns routing behavior with the observed training dynamics. UAR introduces two complementary mechanisms. First, *Evidence-Triggered Expansion* increases the number of active experts for tokens residing in high-energy regions of the Cognitive Vacuum, preserving gradient flow and expanding the search space when evidence is insufficient. Second, *Epistemic Masking* applies load-balancing regularization selectively to high-uncertainty tokens, allowing confident tokens to reinforce expert specialization without being pulled toward global uniformity.

Empirically, we validate UAR across diverse model scales, demonstrating that aligning routing with thermodynamic dynamics significantly improves both convergence stability and representational quality. By dynamically expanding expert capacity via **Evidence-Triggered Routing** to bridge the “Cognitive Vacuum” and applying **Cognitive Load Balancing** to resolve the diversity–specialization dilemma, our framework effectively prevents bias explosion and oscillation. Ultimately, UAR yields superior expert distinctiveness and reduced perplexity with negligible computational overhead ($\sim 1.01 \times$ FLOPs), establishing a robust, physically grounded paradigm for scaling models.

2 Preliminaries

2.1 Mixture-of-Experts

While standard Transformers (Vaswani et al., 2017) couple model capacity with inference cost, Mixture-of-Experts (MoE) decouples them by replacing dense feed-forward blocks with sparse experts $\{\mathcal{E}_k\}_{k=1}^K$. For a token representation h , a router generates gating weights $G(h)$ to activate a sparse subset of experts:

$$\text{MoE}(h) = \sum_{k=1}^K G_k(h) \mathcal{E}_{\theta^{(k)}}(h). \quad (1)$$

Load Balancing. To prevent *expert collapse* (where routing concentrates on few experts), an auxiliary loss is commonly minimized (Lepikhin et al., 2020; Du et al., 2022):

$$\mathcal{L}_{\text{aux}} = N \sum_{k=1}^K f_k P_k, \quad (2)$$

where f_k is the fraction of tokens dispatched to expert k and P_k is the average routing probability. This enforces uniform expert utilization, an assumption we revisit in subsequent sections.

3 Uncertainty in MoE

Routing decisions constitute the primary source of uncertainty in sparse MoE. We formalize this by linking the router’s logit distribution to specific uncertainty modalities.

3.1 Router Entropy: Aleatoric Uncertainty

Let $\mathbf{z} \in \mathbb{R}^K$ be the router logits for token h . The routing probability \mathbf{p} is derived via Softmax with temperature τ , yielding the *Router Entropy*:

$$\mathcal{H}(\mathbf{p}) = - \sum_{k=1}^K p_k \log p_k, \quad p_k = \frac{\exp(z_k/\tau)}{\sum_{j=1}^K \exp(z_j/\tau)} \quad (3)$$

$\mathcal{H}(\mathbf{p})$ quantifies local conflict during expert selection (aleatoric uncertainty). This metric naturally aligns with auxiliary Load Balancing objectives, where higher entropy signifies a diffused routing state desirable for avoiding collapse.

Entropy-Adaptive Routing (Top- p). To dynamically allocate compute, Top- p routing determines the active expert count N based on cumulative probability mass. This creates a functional dependency where N scales implicitly with uncertainty, ranging from $N = 1$ in high-confidence regimes to $N \approx \lceil pK \rceil$ at maximum entropy (see Appendix D for formal derivations). However, these approaches

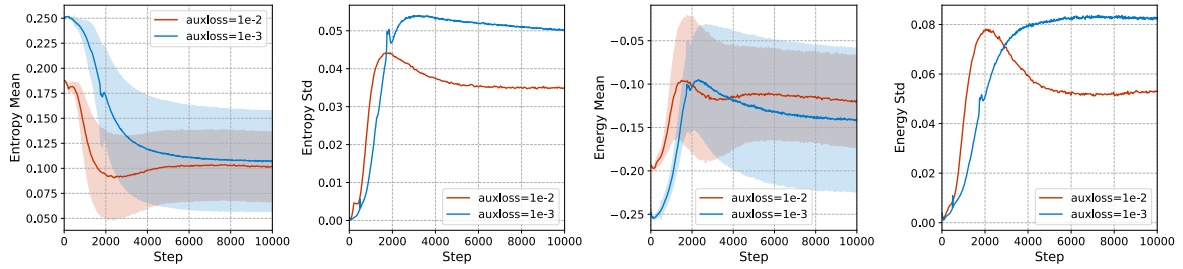


Figure 2: **Entropy and Energy Landscapes.** (Left) Entropy decay illustrates symmetry breaking. (Right) The depth-induced energy gap highlights persistent uncertainty in deeper layers.

158 treat routing as a static classification task, overlooking
 159 the evolutionary nature of routing decisions and
 160 their impact on expert specialization. This lack of
 161 a principled dynamic model prevents the optimal
 162 synchronization of the router learning processes.

163 3.2 The Epistemic Blind Spot

164 While $\mathcal{H}(\mathbf{p})$ effectively measures inter-expert con-
 165 flict, it suffers from a calibration failure regarding
 166 knowledge boundaries. The Softmax operator nor-
 167 malizes *relative* likelihoods, implicitly assuming
 168 the optimal expert exists within the current pool
 169 $\{1, \dots, K\}$.

170 However, for Out-of-Distribution (OOD) tokens
 171 or sparse feature regions, all logits z_k may be small
 172 in magnitude. In such cases, Softmax can still
 173 yield a low-entropy distribution, causing the router
 174 to “confidently” select a suboptimal expert. This
 175 phenomenon masks the *epistemic* uncertainty: the
 176 low-entropy state reflects an absence of internal
 177 competition rather than the presence of strong evi-
 178 dence, ignoring the potential "missing knowledge"
 179 of an ideal, infinite-capacity expert pool (further
 180 detailed in Appendix E).

181 3.3 Routing Energy: Epistemic Uncertainty

182 To capture absolute evidence magnitude, we adopt
 183 a thermodynamic formulation. Defining the *semantic*
 184 *energy* of assigning token h to expert k as $E_k =$
 185 $-z_k$, the routing probability follows the Boltz-
 186 mann distribution $p_k = e^{-E_k/\tau}/Z_\theta$. We quantify
 187 global routing confidence via the *Helmholtz Free*
 188 *Energy* (Helmholtz, 1982):

$$189 \mathcal{F}(\mathbf{z}) = -\tau \log Z_\theta = -\tau \log \sum_{k=1}^K \exp(z_k/\tau) \quad (4)$$

190 Unlike the normalized entropy \mathcal{H} , \mathcal{F} acts as a scalar
 191 proxy for total evidence magnitude. Lower \mathcal{F} indi-
 192 cates high absolute confidence, while higher values

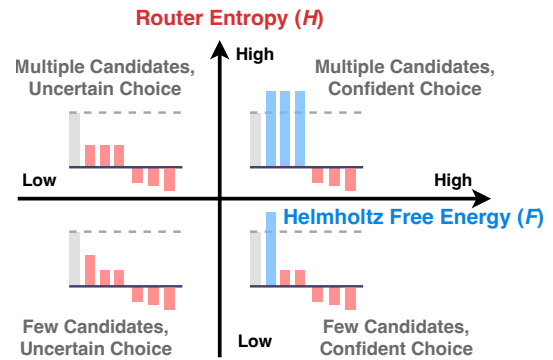


Figure 3: **Thermodynamic Routing Taxonomy.** The \mathcal{H} - \mathcal{F} plane reveals distinct regimes, notably the "Stubborn Bias" (hallucination) and "Functional Ambiguity" (potential collaboration).

193 signal low-evidence regimes (e.g., OOD inputs)
 194 regardless of the distribution’s sharpness.

195 **Taxonomy of Routing States** We characterize
 196 the MoE routing landscape via the interplay be-
 197 tween Entropy (\mathcal{H}) and Free Energy (\mathcal{F}), mapping
 198 token-expert interactions into four distinct quad-
 199 rants (Figure 3): The routing dynamics can be cate-
 200 gorized into four distinct regimes based on the in-
 201 terplay of entropy (\mathcal{H}) and free energy (\mathcal{F}): the ideal
 202 **Specialized Consensus** (low \mathcal{H} , low \mathcal{F}), where
 203 high certainty and strong evidence indicate the to-
 204 ken aligns with a specialized expert domain; **Func-**
 205 **functional Ambiguity** (high \mathcal{H} , low \mathcal{F}), which signals
 206 polysemantic features that possess strong yet con-
 207 flicting evidence, making them suitable for collab-
 208 orative expert activation; **Stubborn Bias** (low \mathcal{H} ,
 209 high \mathcal{F}), a failure mode analogous to hallucination
 210 where the router exhibits unjustified confidence de-
 211 spite lacking absolute evidence, often driven by
 212 overfitting; and **Cognitive Vacuum** (high \mathcal{H} , high
 213 \mathcal{F}), representing a state of maximum uncertainty
 214 devoid of both preference and evidence, a condition
 215 prevalent during initialization or in deeper layers.

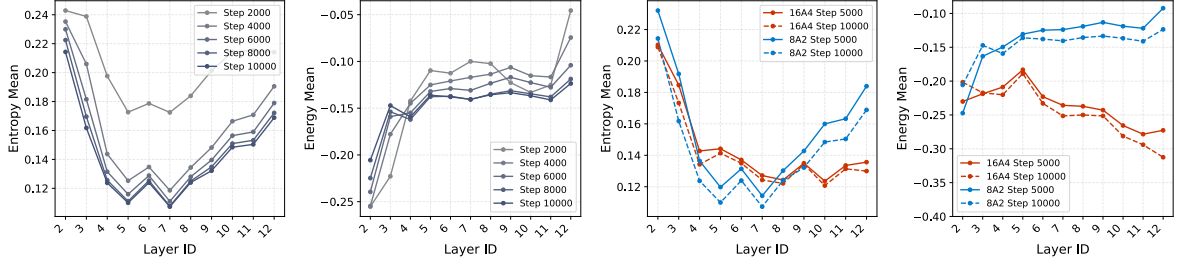


Figure 4: **Left panels:** temporal dynamics of Router Entropy and Energy Mean from step 2000 to 10000. **Right panels:** The impact of expert granularity (Fine-grained 16A4 vs. Coarse 8A2 in same parameter setting).

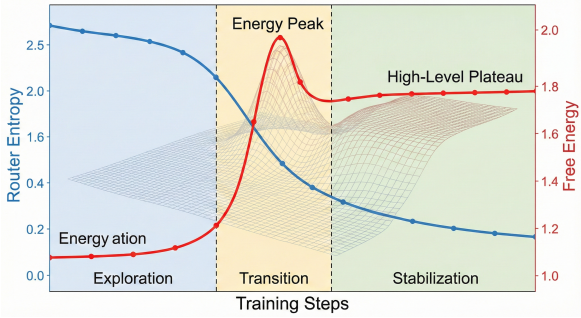


Figure 5: **Thermodynamic Trajectory.** Training evolves through three phases: (A) Exploration, (B) Transition (marked by a sharp **Energy Peak**), and (C) Stabilization into a **High-Level Plateau**, indicating persistent epistemic uncertainty despite entropy minimization.

4 Observation of Training Dynamics

We employ two 1.2B total parameter MoE Transformers (25% activation, 8A2, 16A4) trained for 50B data. Auxiliary loss coefficients $\lambda \in \{10^{-2}, 10^{-3}\}$ are varied to probe the impact of regularization on the epistemic landscape.

4.1 Symmetry Breaking and Transitions

Symmetry Breaking. Mean routing entropy $\bar{\mathcal{H}}$ in Figure 2 decay illustrates the router’s transition from initialization (State IV) to specialization. A peak in $\text{std}(\mathcal{H})$ (0k–2k steps) marks a *critical phase transition* distinguishing "easy" tokens from "hard" ones. While shallow layers converge to State I (Specialized Consensus), deep layers maintain a high entropy floor in Figure 4, indicating a persistent regime of functional ambiguity.

Depth-Induced Energy Gap. While global $\bar{\mathcal{F}}$ decreases with epistemic maturation, deep layers retain significantly higher energy and variance in Figure 4. This "depth-dependent chaos" confirms a *Cognitive Vacuum*: routers in deep layers distribute tokens (low \mathcal{H}) without absolute evidence (high \mathcal{F}).

This discrepancy necessitates dynamic adjustment to compensate for the lack of certainty.

Diversity-Specialization Dilemma. A fine-grained analysis reveals a trade-off between utilization balance and representation distinctiveness (Figure 2). High λ enforces uniformity at the cost of expert cosine similarity, yielding redundant generalists, whereas low λ risks capacity under-utilization.

Thermodynamic Evolution and Capacity Gap.

Figure 4 reveals a distinct U-shaped entropy profile, where middle layers achieve specialized consensus (State I) before regressing to uncertainty in deeper layers. Crucially, this trajectory exposes a capacity-driven divergence: while the fine-grained 16A4 consistently minimizes Helmholtz Free Energy \mathcal{F} , the coarse-grained 8A2 suffers from a sharper entropy rebound coupled with *Energy Saturation*. This simultaneous elevation of entropy and energy confirms that 8A2 collapses into a *Cognitive Vacuum* (State IV) in deep layers—lacking the structural capacity to resolve semantic ambiguity—whereas 16A4 provides sufficient expressivity to mitigate this epistemic risk.

4.2 Thermodynamic Evolution

We track the thermodynamic trajectory of the router across three distinct phases (Figure 8). In **Phase A (Initialization)**, \mathcal{F} remains in a quasistationary low state alongside maximum \mathcal{H} , reflecting disordered potentials. Upon entering **Phase B (Transition)**, \mathcal{F} exhibits a non-monotonic *Energy Peak*, rising sharply as the system attempts to break symmetry. Finally, in **Phase C (Convergence)**, \mathcal{F} settles into a *High-Level Plateau* significantly above its initial value, particularly in deeper layers.

Frustrated Exploration. The Energy Peak signifies a regime of structural tension. Initially, weight

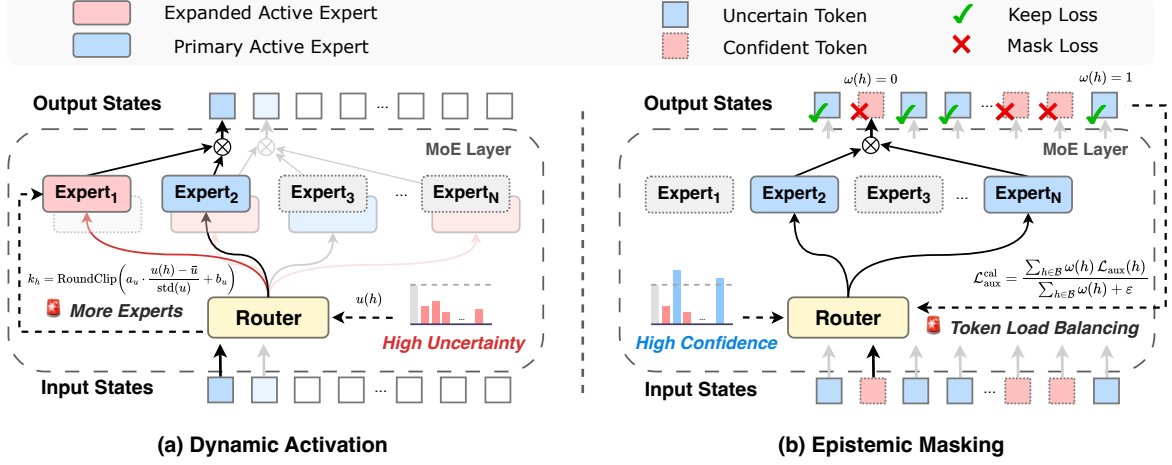


Figure 6: **Uncertainty-aware MoE routing framework.** (a) **Dynamic Activation:** The expert budget k_h is adaptively scaled by token uncertainty $u(h)$ to form a robust ensemble. (b) **Epistemic Masking:** Load-balancing loss is selectively applied via $\omega(h)$ to uncertain tokens, decoupling expert exploration from specialization.

decay constrains logits near zero ($\sum e^{z_k/\tau} \approx K$). As specialization begins, a conflict emerges: task gradients attempt to amplify specific expert weights, while the auxiliary loss \mathcal{L}_{aux} pulls them toward uniformity. This "tug-of-war" creates a thermodynamic barrier—the router is stripped of initial uniform evidence but prevented from establishing high-confidence basins.

The Cognitive Vacuum. The terminal High-Level Plateau exposes a critical epistemic failure. Even at convergence, the router minimizes relative conflict ($\mathcal{H} \rightarrow 0$) without minimizing absolute free energy. This creates a deceptive state where the model learns to *avoid* clearly incorrect experts rather than *affirm* correct ones—routing with apparent certainty but hollow evidence.

5 Uncertainty-Aware Routing (UAR)

The persistence of the Energy Peak and the depth-dependent plateau demonstrates that optimal search breadth is state-dependent. To mitigate this vacuum, we propose **Uncertainty-Aware Routing (UAR)**, which explicitly utilizes Uncertainty as a control signal to orchestrate a synergistic co-evolution between routing decisions and dynamic expert capacity. UAR modulates search breadth and regularization pressure dynamically, utilizing the router’s internal thermodynamic state to orchestrate the exploration-exploitation trade-off.

5.1 Adaptive Expert Expansion (AE)

Evidence-Triggered Computation. Standard routing assumes constant information density, forcing sparse activation even in low-evidence regimes (e.g., Phase B or State IV). This rigidity degrades gradient signals and induces hallucinations. To "hedge" against epistemic risk, we dynamically scale the active expert count k_h based on token uncertainty $u(h)$ (instantiated as \mathcal{F} or \mathcal{H}). We formalize the activation budget via a normalized affine transformation:

$$k_h = \text{RoundClip}\left(\alpha \cdot \frac{u(h) - \bar{u}}{\sigma_u} + \beta\right), \quad (5)$$

$$\text{RoundClip}(x) = \text{clamp}(\lfloor x \rfloor, K_{\min}, K_{\max}).$$

where \bar{u} and σ_u are moving statistics of $u(h)$. This mechanism effectively converts the router into a *soft, evidence-weighted ensemble* specifically when the signal-to-noise ratio is low. During the critical Transition Phase (Phase B), the resulting spike in k_h prevents premature collapse; in deeper layers (State IV), it acts as a multi-sample estimator to stabilize representation variance.

5.2 Calibrated Load Balancing (CLB)

Decoupling Exploration from Specialization. Standard auxiliary losses \mathcal{L}_{aux} penalize imbalance uniformly, causing *structural interference* where global load constraints degrade local expert specialization (State I). We introduce **Epistemic Masking** to enforce load balancing *selectively*.

We define an epistemic gating weight $\omega(h) = g(u(h), \theta_t)$, where g is a sigmoid shaping function

Model-Params	FLOPs	Commonsense & Reading Comprehension						Continued		LM	Knowledge	Avg.
		SciQ	PIQA	WG	ARC-E	ARC-C	Hella.	LogiQA	BoolQ	Lam.	MMLU	
MoE-8A2 UAR-AE vs. Router Variants (with Auxloss)												
Top- k	1×	52.4	63.2	51.6	42.0	22.4	31.1	21.0	54.6	37.1	25.4	41.7
Top- p	0.98×	53.5	63.0	50.9	41.8	24.1	30.4	23.0	56.0	36.7	25.6	42.2
ReLU router	1.1×	53.7	63.8	50.8	41.0	23.2	32.0	21.2	55.0	38.3	26.1	42.1
UAR-AE	1.06×	51.5	63.6	52.2	40.6	22.7	30.5	24.0	59.0	39.1	26.8	42.6
MoE-8A2 UAR-CLB: Auxloss / Auxfree × CLB (Top- k)												
Auxloss	1×	52.4	63.2	51.6	42.0	22.4	31.1	21.0	54.6	37.1	25.4	41.7
Auxloss + CLB	1×	53.4	63.9	51.9	41.0	23.4	31.0	23.5	60.6	39.1	26.2	43.1
Auxfree	1×	52.3	63.6	50.0	40.6	22.7	30.8	20.3	59.2	37.8	25.8	41.9
Auxfree + CLB	1×	53.0	64.3	52.3	41.9	24.1	31.1	21.2	59.2	38.3	26.5	42.8
MoE-256A8 UAR-AE vs. Router Variants (with Auxloss)												
Top- k	1×	53.9	64.0	51.4	40.4	20.9	30.7	20.9	59.4	38.4	26.0	42.2
Top- p	0.97×	54.1	63.8	51.2	41.0	21.5	30.5	22.1	59.8	38.0	26.2	42.4
ReLU router	1.09×	53.5	64.1	50.9	40.2	21.0	31.5	21.5	58.5	39.0	26.5	42.3
UAR-AE	1.08×	54.8	64.5	52.5	41.2	21.2	31.2	23.5	60.2	39.5	27.0	43.2
MoE-256A8 UAR-CLB: Auxloss / Auxfree × CLB (Top- k)												
Auxloss	1×	53.9	64.0	51.4	40.4	20.9	30.7	20.9	59.4	38.4	26.0	42.2
Auxloss + CLB	1×	54.5	64.8	52.1	41.5	21.8	31.3	22.4	60.5	39.8	26.9	43.1
Auxfree	1×	54.0	63.9	50.8	40.5	21.0	30.9	20.5	59.5	38.6	26.1	42.1
Auxfree + CLB	1×	55.1	65.0	52.9	42.1	22.2	31.8	22.0	60.8	40.1	27.2	43.5

Table 1: **Main Results on Downstream Benchmarks.** We report zero-shot and few-shot accuracy across 10 datasets for MoE-8A2 and MoE-256A8 models. We directly compare UAR with other dynamic routing methods, while CLB is applied as an attachment to existing load-balancing optimization strategies.

centered at a time-annealed threshold θ_t . The calibrated objective focuses regularization solely on uncertain tokens:

$$\mathcal{L}_{\text{aux}}^{\text{cal}} = \frac{\sum_{h \in \mathcal{B}} \omega(h) \mathcal{L}_{\text{aux}}(h)}{\sum_{h \in \mathcal{B}} \omega(h) + \epsilon} \quad (6)$$

Complementary to this, we employ an **Uncertainty-Weighted Bias Update**. Expert bias terms b_k are updated with a distinct learning rate schedule $\lambda(h)$ dependent on $u(h)$:

$$\Delta b_k \propto [\lambda_{\text{hi}} \mathbb{1}(u > \theta_t) + \lambda_{\text{lo}} \mathbb{1}(u \leq \theta_t)] (\text{load}_{\text{tgt}} - \text{load}_k) \quad (7)$$

with $\lambda_{\text{hi}} \gg \lambda_{\text{lo}}$. This ensures that high-uncertainty tokens (States II/IV) drive exploration and load equalization, while confident tokens (State I) are shielded from balancing pressure. This separation preserves the sharpness of learned specializations while preventing pathological collapse in ambiguous regions of the data manifold.

6 Experiments

6.1 Setup

Data. To pretrain UAR models and baseline models, we use OLMo 2 Mix 1124, which matches the OLMoE (Muennighoff et al., 2024) pretraining set. It combines DCLM (Li et al., 2024), Dolma 1.7 (Soldaini et al., 2024) subsets (arXiv, OpenWebMath, Algebraic Stack, peS2o, Wikipedia), and

StarCoder (Li et al., 2023). The dataset totals 3.4 trillion training tokens and we sample 2 million tokens for validation.

Training. All models are trained using Megatron on NVIDIA H800 GPUs, with a sequence length of 4096 and a global batch size of 256. We adopt a cosine learning rate decay schedule combined with a multi-step learning rate scheduler, starting from a peak learning rate of 5×10^{-4} and decaying to 10% of the initial value, with a warmup over the first 5% of total tokens. Model weights are initialized with standard deviation $\sqrt{2/(5d)}$, where d is the hidden size. Mixture-of-experts models use an auxiliary loss coefficient of 10^{-3} . Each model is trained with a data-to-parameter ratio of 200, on the same data volume, from random initialization.

Evaluation. We employed the lm-evaluation-harness (Gao et al., 2023) to evaluate our models. For common sense and reading comprehension tasks, we report 0-shot accuracy for SciQ (Welbl et al., 2017), PIQA (Bisk et al., 2020), WinoGrande (WG) (Sakaguchi et al., 2020), ARC Easy (ARC-E) (Clark et al., 2018a), and 10-shot HellaSwag (Hella.) (Zellers et al., 2019), alongside 25-shot accuracy for ARC Challenge (ARC-C) (Clark et al., 2018b). For continued QA and text understanding, we report 0-shot accuracy for LogiQA (Liu et al., 2020), 32-shot BoolQ (Clark et al., 2019), and 0-

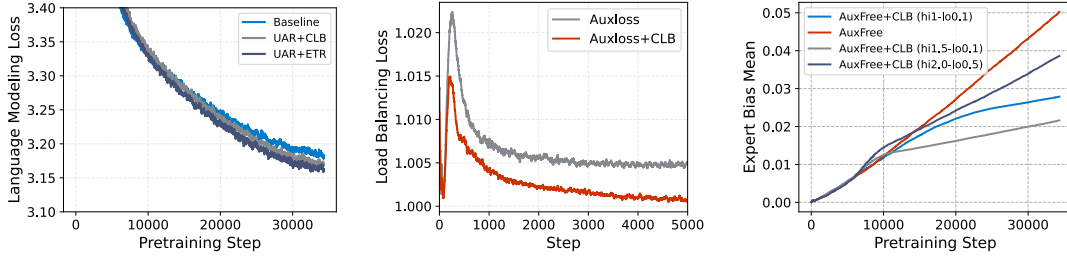


Figure 7: **Performance Gains and Training Stability with UAR.** **Left:** Language modeling loss curves. **Middle:** Comparison of load balancing loss trajectories. **Right:** Evolution of expert bias mean in AuxFree settings, .

shot LAMBADA (Lam.) (Paperno et al., 2016). All reported results are calculated with the mean and stderr of multiple experiments.

Baseline. Our primary analysis is conducted on a 0.18B-parameter model, with additional results reported at the 0.5B scale. All models adopt the DeepSeekMoE backbone (Dai et al., 2024) without expert sharing: the first-layer FFN remains dense, while all subsequent FFNs are replaced by MoE modules with a Softmax router. We compare our dynamic routing method against Top- k , Top- p , and ReMoE (Wang et al., 2024b). The proposed Epistemic Masking (EM) is evaluated under both loss-free and auxiliary loss-based load-balancing regimes to analyze routing behavior across different stabilization strategies.

6.2 Result

Performance and Load Balancing. Table 1 demonstrates that Uncertainty-Aware Evidence-Triggered Routing (UAR-AE) consistently outperforms fixed Top- k strategies and matches dynamic baselines (e.g., Top- p , ReLU) under aligned FLOPs. By leveraging entropy to adaptively allocate computational budget, UAR-AE enhances representational quality without relying on explicit auxiliary losses. Furthermore, our CLB serves as a zero-cost, parameter-free refinement. Whether applied to standard auxiliary-loss (Auxloss) or auxiliary-free (Auxfree) regimes, CLB universally improves downstream accuracy and convergence.

Scalability. We assess robustness across active parameters ($N \in [182\text{M}, 975\text{M}]$) and expert counts ($E \in [4, 128]$). As shown in Figure 7, UAR-based architectures consistently surpass standard MoE baselines. Notably, this performance advantage amplifies in larger model regimes and persists as the expert pool expands, indicating that UAR effectively mitigates the saturation and redundancy

Method	Perplexity ↓	MaxVio ↓
<i>Setting A: With Auxiliary Loss</i>		
Auxloss (Baseline)	23.93	0.152
Uncertainty = Energy		
Start@10k, Warmup=5k	23.77 (-0.16)	0.108 (-0.044)
Start@20k, Warmup=5k	23.79 (-0.14)	0.112 (-0.040)
Uncertainty = Entropy		
Start@10k, Warmup=10k	23.75 (-0.18)	0.105 (-0.047)
Start@20k, Warmup=10k	23.78 (-0.15)	0.109 (-0.043)
<i>Setting B: AuxFree</i>		
AuxFree (Baseline)	23.85	0.102
Uncertainty = Energy		
Start@10k, Warmup=5k	23.69 (-0.16)	0.074 (-0.028)
Uncertainty = Entropy		
Start@10k, Warmup=10k	23.67 (-0.18)	0.069 (-0.033)

Table 2: Ablation for Calibrated Load Balancing on MoE-162M. We compare two settings: standard **Auxloss** and **AuxFree** (bias-only).

often observed in massive expert populations.

Stability and Thermodynamics. Experiments on 1B and 3B models confirm that CLB significantly reduces global load violation ($\text{MaxVio}_{\text{global}}$) and validation perplexity. As detailed in Appendix G, this approach prevents the characteristic "early oscillation / late collapse" pattern of MoE training. Moreover, analyzing the thermodynamic landscape reveals that UAR fundamentally stabilizes routing dynamics. By lowering Helmholtz free energy \mathcal{F} and smoothing routing entropy \mathcal{H} , the framework avoids high-uncertainty "Cognitive Vacuum" states, ensuring expert selection aligns strictly with the model's epistemic confidence.

6.3 Ablation Studies

Ablation: AE Efficiency and Proxy Selection. Table 3 validates that AE yields consistent perplexity reductions with marginal computational overhead. Specifically, Energy-based AE ($a = 1.0$) surpasses the fixed Top-2 baseline by 0.19 PPL while incurring only a $1.03\times$ increase in FLOPs; a

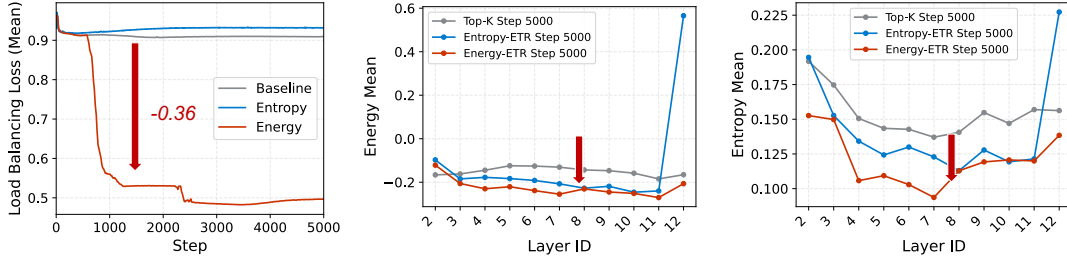


Figure 8: **Uncertainty Metrics on Training Dynamics.** **Left:** The evolution of load balancing loss throughout training. **Middle:** Layer-wise distribution of Energy Mean at step 5000. **Right:** Layer-wise distribution of Entropy Mean at step 5000.

conservative setting ($a = 0.5$) achieves a 0.14 PPL gain with negligible cost ($1.01 \times$ FLOPs). Crucially, Energy (\mathcal{F}) consistently outperforms Entropy (\mathcal{H}) as an uncertainty proxy (e.g., Δ PPL -0.14 vs. -0.11 at $a = 0.5$), confirming that \mathcal{F} more accurately captures the absolute epistemic uncertainty required to target "hard" tokens in the Cognitive Vacuum.

Ablation: CLB Robustness. Table 2 assesses CLB across standard auxiliary-loss (Setting A) and bias-only (Setting B, AuxFree) regimes. In Setting A, CLB significantly reduces MaxVio by ~ 0.044 while improving perplexity ($23.93 \rightarrow 23.77$), demonstrating that masking confident tokens minimizes interference with expert specialization. The benefits amplify in the AuxFree regime, where CLB achieves the lowest observed violation ($\text{MaxVio} \approx 0.062$) and superior perplexity (23.69). This consistent reduction across settings proves CLB enforces load uniformity effectively without compromising convergence.

6.4 Analysis

We now dissect the mechanisms behind UAR's performance, linking empirical gains to the thermodynamic dynamics discussed in Sec. 4. Detailed qualitative analysis is provided in Appendix H.

Resolving the Cognitive Vacuum via Free Energy. Fig. 8 corroborates our hypothesis that Router Entropy fails to capture the "Cognitive Vacuum" in deep layers (8–12). While entropy-based baselines struggle with persistent ambiguity, UAR utilizes Free Energy (\mathcal{F}) as a control signal to explicitly target these high-energy regions. By dynamically expanding expert slots, UAR suppresses the energy floor (Fig. 8, Middle), converting epistemic uncertainty into robust multi-expert ensembles rather than forcing low-confidence routing.

Buffering Phase Transitions. Fig. 9 validates that UAR acts as an adaptive buffer during the "Frustrated Exploration" phase (steps 800–2000). The effective number of experts (k_{eff}) automatically peaks at ≈ 2.1 during this "Symmetry Breaking" window, aligning with periods of high thermodynamic stress. This confirms that the model naturally hedges optimization risks by ensembling experts when the landscape is chaotic, before decaying toward sparsity in the stabilization phase without manual scheduling.

Harmonizing Specialization and Uniformity. Finally, Fig. 7 demonstrates how CLB resolves the "Diversity-Specialization Dilemma." Unlike the AuxFree baseline which suffers from "Bias Explosion" (Red line), CLB's Epistemic Masking suppresses structural interference. By regularizing only high-uncertainty tokens, the system allows confident states to reinforce specialization while leveraging uncertain states for load balancing. This selective pressure yields a cleaner optimization path and lower convergence loss, achieving uniformity without impeding expert maturation.

7 Conclusion

Standard MoE training is hindered by "Frustrated Exploration" and persistent uncertainty in deep layers, caused by misalignment between static routing and evolving model confidence. To address this, we introduce Uncertainty-Aware Routing, which adaptively expands expert activation and selectively applies load balancing based on the router's thermodynamic state. Our results demonstrate that dynamically aligning routing decisions with model uncertainty improves both performance and expert specialization, providing a principled path for efficient and robust MoE scaling.

8 Limitations

While UAR demonstrates improved performance and stability in MoE training, several limitations remain. First, the approach relies on appropriate calibration of uncertainty metrics (such as Helmholtz Free Energy), which may require additional tuning across different architectures and datasets. Second, the dynamic adjustment of expert activation introduces extra computation and memory overhead, especially when the number of active experts becomes large for uncertain tokens. Third, our evaluation is primarily conducted on language modeling tasks; the generalizability of UAR to other modalities or tasks (e.g., vision, multi-modal models, or reinforcement learning) requires further study. Finally, although UAR reduces the tradeoff between load balancing and specialization, achieving the optimal balance remains challenging in extremely large-scale or highly imbalanced expert settings. This manuscript was prepared with the assistance of AI tools for language polishing and structural refinement; all scientific content and conclusions are solely the responsibility of the authors.

References

- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450*.
- Henning Bartsch, Ole Jorgensen, Domenic Rosati, Jason Hoelscher-Obermaier, and Jacob Pfau. 2023. Self-consistency of large language models under ambiguity. *arXiv preprint arXiv:2310.13439*.
- Christopher M Bishop and Nasser M Nasrabadi. 2006. *Pattern recognition and machine learning*, volume 4. Springer.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2020. [PIQA: reasoning about physical commonsense in natural language](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 7432–7439. AAAI Press.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

- Yanxi Chen, Xuchen Pan, Yaliang Li, Bolin Ding, and Jingren Zhou. 2024. [Ee-llm: Large-scale training and inference of early-exit large language models with 3d parallelism](#). 566
567
568
569
- Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. 2019. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*. 570
571
572
573
574
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018a. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*. 575
576
577
578
579
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018b. [Think you have solved question answering? try arc, the AI2 reasoning challenge](#). *CoRR*, abs/1803.05457. 580
581
582
583
584
- Damai Dai, Chengqi Deng, Chenggang Zhao, RX Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng, Xingkai Yu, Yu Wu, et al. 2024. Deepseek-moe: Towards ultimate expert specialization in mixture-of-experts language models. *arXiv preprint arXiv:2401.06066*. 585
586
587
588
589
590
- Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. 2022. Glam: Efficient scaling of language models with mixture-of-experts. In *International conference on machine learning*, pages 5547–5569. PMLR. 591
592
593
594
595
596
- Mostafa Elhoushi, Akshat Shrivastava, Diana Liskovich, Basil Hosmer, Bram Wasti, Liangzhen Lai, Anas Mahmoud, Bilge Acun, Saurabh Agarwal, Ahmed Roman, Ahmed Aly, Beidi Chen, and Carole-Jean Wu. 2024. [Layerskip: Enabling early exit inference and self-speculative decoding](#). page 12622–12642. 597
598
599
600
601
602
- William Fedus, Barret Zoph, and Noam Shazeer. 2022. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *The Journal of Machine Learning Research*, 23(1):5232–5270. 603
604
605
606
607
- Yuchen Feng, Bowen Shen, Naibin Gu, Jiaxuan Zhao, Peng Fu, Zheng Lin, and Weiping Wang. 2025. Dive into moe: Diversity-enhanced reconstruction of large language models from dense into mixture-of-experts. *arXiv preprint arXiv:2506.09351*. 608
609
610
611
612
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. 2023. [A framework for few-shot language model evaluation](#). 613
614
615
616
617
618
619
620
621

622	Xiang Gao, Jiaxin Zhang, Lalla Mouatadid, and Kamalika Das. 2024. Spuq: Perturbation-based uncertainty quantification for large language models. <i>arXiv preprint arXiv:2403.02509</i> .	676
623		677
624		678
625		679
		680
626	Yashvir S Grewal, Edwin V Bonilla, and Thang D Bui. 2024. Improving uncertainty quantification in large language models via semantic embeddings. <i>arXiv preprint arXiv:2410.22685</i> .	681
627		682
628		683
629		684
630	Hermann Helmholtz. 1982. <i>Über die Erhaltung der Kraft</i> , volume 1. Walter de Gruyter GmbH & Co KG.	685
631		686
632		687
633	Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. 2022. Training compute-optimal large language models. <i>arXiv preprint arXiv:2203.15556</i> .	688
634		689
635		690
636		691
637		692
638		693
639	Quzhe Huang, Zhenwei An, Nan Zhuang, Mingxu Tao, Chen Zhang, Yang Jin, Kun Xu, Liwei Chen, Songfang Huang, and Yansong Feng. 2024. Harder tasks need more experts: Dynamic routing in moe models. <i>arXiv preprint arXiv:2403.07652</i> .	694
640		695
641		696
642		697
643		698
644	A Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Benoît Savary, Charles Bamford, Devendra Singh Chaplot, Daniele de la Casas, Emily Bressand Hanna, François Bressand, et al. 2024. Mixtral of experts. <i>arXiv preprint arXiv:2401.04088</i> .	699
645		700
646		701
647		702
648		703
649		704
650	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. <i>arXiv preprint arXiv:2001.08361</i> .	705
651		706
652		707
653		708
654		709
655	Elizaveta Kostenok, Daniil Cherniavskii, and Alexey Zaytsev. 2023. Uncertainty estimation of transformers’ predictions via topological analysis of the attention matrices. <i>arXiv preprint arXiv:2308.11295</i> .	710
656		711
657		712
658		713
659	Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2024. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. <i>Nature</i> .	714
660		715
661		716
662		717
663	Dmitry Lepikhin, Hyoungho Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. 2020. Gshard: Scaling giant models with conditional computation and automatic sharding. <i>arXiv preprint arXiv:2006.16668</i> .	718
664		719
665		720
666		721
667		722
668		723
669	Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Yitzhak Gadre, Hritik Bansal, Etash Guha, Sedrick Scott Keh, Kushal Arora, et al. 2024. Datacomp-lm: In search of the next generation of training sets for language models. <i>Advances in Neural Information Processing Systems</i> , 37:14200–14282.	724
670		725
671		726
672		727
673		728
674		729
675		730
	Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. 2023. Starcoder: may the source be with you! <i>arXiv preprint arXiv:2305.06161</i> .	
	Yinghao Li, Rushi Qiang, Lama Moukheiber, and Chao Zhang. 2025. Language model uncertainty quantification with attention chain. <i>arXiv preprint arXiv:2503.19168</i> .	
	Zhen Lin, Shubhendu Trivedi, and Jimeng Sun. 2023. Generating with confidence: Uncertainty quantification for black-box large language models. <i>arXiv preprint arXiv:2305.19187</i> .	
	Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2020. Logiqa: A challenge dataset for machine reading comprehension with logical reasoning. <i>arXiv preprint arXiv:2007.08124</i> .	
	Linyu Liu, Yu Pan, Xiaocheng Li, and Guanting Chen. 2024. Uncertainty estimation and quantification for llms: A simple supervised approach, 2024. <i>URL https://arxiv.org/abs/2404.15993</i> .	
	Zichang Liu, Jue Wang, Tri Dao, Tianyi Zhou, Binhang Yuan, Zhao Song, Anshumali Shrivastava, Ce Zhang, Yuandong Tian, Christopher Re, and Beidi Chen. 2023. <i>Deja vu: Contextual sparsity for efficient llms at inference time</i> . <i>arXiv preprint arXiv:2310.17157</i> .	
	Qing Lyu, Kumar Shridhar, Chaitanya Malaviya, Li Zhang, Yanai Elazar, Niket Tandon, Marianna Apidianaki, Mrinmaya Sachan, and Chris Callison-Burch. 2025. Calibrating large language models with sample consistency. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 19260–19268.	
	Yongyu Mu, Yuzhang Wu, Yuchun Fan, Chenglong Wang, Hengyu Li, Qiaozhi He, Murun Yang, Tong Xiao, and Jingbo Zhu. 2024. <i>Cross-layer attention sharing for large language models</i> . <i>Preprint, arXiv:2408.01890</i> .	
	Niklas Muennighoff, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Jacob Morrison, Sewon Min, Weijia Shi, Pete Walsh, Oyvind Tafjord, Nathan Lambert, et al. 2024. Olmoe: Open mixture-of-experts language models. <i>arXiv preprint arXiv:2409.02060</i> .	
	OpenAI. 2023. Gpt-4 technical report. <i>ArXiv</i> , page abs/2303.08774.	
	Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. 2016. The lambada dataset: Word prediction requiring a broad discourse context. <i>arXiv preprint arXiv:1606.06031</i> .	
	Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Mario Neumann, Rodolphe Jenatton, António Susano Pinto, Daniel Keysers, and Neil Houlsby. 2021. Scaling	

731	vision with sparse mixture of experts. In <i>Advances in Neural Information Processing Systems</i> , volume 34, pages 8583–8595.	
732		
733		
734	Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. Winogrande: An adversarial winograd schema challenge at scale . In <i>The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020</i> , pages 8732–8740. AAAI Press.	
735		
736		
737		
738		
739		
740		
741		
742		
743		
744	Claude E Shannon. 1948. A mathematical theory of communication. <i>The Bell system technical journal</i> , 27(3):379–423.	
745		
746		
747	Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. 2017. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. <i>arXiv preprint arXiv:1701.06538</i> .	
748		
749		
750		
751		
752	Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, et al. 2024. Dolma: An open corpus of three trillion tokens for language model pretraining research. In <i>Proceedings of the 62nd annual meeting of the association for computational linguistics (volume 1: long papers)</i> , pages 15725–15788.	
753		
754		
755		
756		
757		
758		
759		
760	Linwei Tao, Yi-Fan Yeh, Minjing Dong, Tao Huang, Philip Torr, and Chang Xu. 2025. Revisiting uncertainty estimation and calibration of large language models. <i>arXiv preprint arXiv:2505.23854</i> .	
761		
762		
763		
764	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	
765		
766		
767		
768		
769	Jingcun Wang, Yu-Guang Chen, Ing-Chao Lin, Bing Li, and Grace Li Zhang. 2024a. Basis sharing: Cross-layer parameter sharing for large language model compression . <i>Preprint</i> , arXiv:2410.03765.	
770		
771		
772		
773	Ziteng Wang, Jun Zhu, and Jianfei Chen. 2024b. Re-moe: Fully differentiable mixture-of-experts with relu routing. <i>arXiv preprint arXiv:2412.14711</i> .	
774		
775		
776	Tianwen Wei, Bo Zhu, Liang Zhao, Cheng Cheng, Biye Li, Weiwei Lü, Peng Cheng, Jianhao Zhang, Xiaoyu Zhang, Liang Zeng, et al. 2024. Skywork-moe: A deep dive into training techniques for mixture-of-experts language models. <i>arXiv preprint arXiv:2406.06563</i> .	
777		
778		
779		
780		
781		
782	Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017. Crowdsourcing multiple choice science questions . In <i>Proceedings of the 3rd Workshop on Noisy User-generated Text, NUT@EMNLP 2017, Copenhagen, Denmark, September 7, 2017</i> , pages 94–106. Association for Computational Linguistics.	
783		
784		
785		
786		
787		
	Quan Xiao, Debarun Bhattacharjya, Balaji Ganesan, Radu Marinescu, Katsiaryna Mirylenka, Nhan H Pham, Michael Glass, and Junkyu Lee. 2025. The consistency hypothesis in uncertainty quantification for large language models. <i>arXiv preprint arXiv:2506.21849</i> .	788 789 790 791 792 793
	Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. 2023. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. <i>arXiv preprint arXiv:2306.13063</i> .	794 795 796 797 798
	F. Xue, Z. Zheng, Y. Fu, J. Ni, and W. Zhou. 2024. Openmoe: An early effort on open mixture-of-experts language models . <i>arXiv preprint arXiv:2402.01739</i> .	799 800 801 802
	An Yang, Anpeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. 2025. Qwen3 technical report. <i>arXiv preprint arXiv:2505.09388</i> .	803 804 805 806
	Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? In <i>Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers</i> , pages 4791–4800. Association for Computational Linguistics.	807 808 809 810 811 812 813 814
	Jun Zhang, Desen Meng, Ji Qi, Zhenpeng Huang, Tao Wu, and Limin Wang. 2024. p-mod: Building mixture-of-depths mllms via progressive ratio decay . <i>Preprint</i> , arXiv:2412.04449.	815 816 817 818
	Tunyu Zhang, Haizhou Shi, Yibin Wang, Hengyi Wang, Xiaoxiao He, Zhuowei Li, Haoxian Chen, Ligong Han, Kai Xu, Huan Zhang, et al. 2025. Token-level uncertainty estimation for large language model reasoning. <i>arXiv preprint arXiv:2505.11737</i> .	819 820 821 822 823
	Yutian Zhou, Tao Lei, Henry Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew M Dai, Quoc V Le, James Laudon, et al. 2022. Mixture-of-experts with expert choice routing. In <i>Advances in Neural Information Processing Systems</i> .	824 825 826 827 828

A Appendix

A Preliminaries

A.1 Transformer

A Transformer (Vaswani et al., 2017) consists of L stacked layers with parameters $\Theta = \{\theta_1, \dots, \theta_L\}$. Given an input sequence $\mathbf{x} \in \mathbb{R}^{N \times d}$, each layer applies multi-head self-attention followed by a feed-forward transformation under layer normalization (Ba et al., 2016):

$$H' = \text{MHSA}(\text{LN}(H)). \quad (8)$$

A standard Transformer computes its output through sequential composition,

$$\mathcal{M}_{\text{std}}(\mathbf{x}) = \psi \circ \mathcal{T}_{\theta_L} \circ \dots \circ \mathcal{T}_{\theta_1} \circ \phi(\mathbf{x}), \quad (9)$$

where ϕ and ψ denote the embedding and output heads. Both parameter count and computational cost scale linearly with depth, tightly coupling model capacity and inference cost.

A.2 Mixture-of-Experts

Mixture-of-Experts layers replace dense feed-forward blocks with sparse conditional computation. For a token representation h , a router maps it to expert weights over K experts, producing

$$\text{MoE}(h) = \sum_{k=1}^K G_k(h) \mathcal{E}_{\theta^{(k)}}(h), \quad (10)$$

where only a small subset of experts is activated per token. This design decouples total parameter count from per-token computation.

Load Balancing. A central challenge in MoE training is expert collapse, where routing concentrates on a few experts. To prevent this, an auxiliary load-balancing loss is commonly introduced (Lepikhin et al., 2020; Du et al., 2022). For a batch of B tokens, let f_k denote the fraction of tokens routed to expert k , and P_k the average routing probability. The auxiliary objective

$$\mathcal{L}_{\text{aux}} = N \sum_{k=1}^K f_k P_k \quad (11)$$

encourages uniform expert utilization by penalizing deviations from balanced routing. Minimizing this loss drives both f_k and P_k toward uniformity, implicitly maximizing routing entropy. While effective at preventing collapse, this formulation assumes that uniformity is desirable throughout training, an assumption we revisit in subsequent sections. Extended definitions in Appendix C.

B Related Work

Uncertainty estimation methods. Recently, numerous uncertainty estimation methods for LLMs have been proposed. These include methods that utilize natural language for uncertainty feedback, including heuristically designed and trained approaches (Tao et al., 2025; Xiong et al., 2023; Lin et al., 2023); methods that estimate uncertainty based on model states, including those leveraging prior knowledge or statistical observations of model states (Kostenok et al., 2023; Li et al., 2025; Liu et al., 2024), or observing changes under perturbations (Zhang et al., 2025; Gao et al., 2024); and methods that take into account the semantics of the response, including consistency-based uncertainty characterizations (Lyu et al., 2025; Bartsch et al., 2023; Xiao et al., 2025) and approaches that integrate semantics with model states (Kuhn et al., 2024; Grewal et al., 2024).

Dynamic Computation Allocation Dynamic Computation Allocation, like Mixture-of-Expert (MoE), reduce computational overhead by activating only a subset of networks (Fedus et al., 2022; Riquelme et al., 2021; Zhou et al., 2022; Jiang et al., 2024; Xue et al., 2024). Some works focus on elastic computation in depth, such as early exit (Elhoushi et al., 2024; Chen et al., 2024), parameter sharing (Mu et al., 2024; Wang et al., 2024a) or using token-routing for dynamic layer skipping (Zhang et al., 2024).

C Extended Preliminaries

C.1 Transformer Architecture

A Transformer is a deep neural network model for sequence modeling, defined by a sequence of L stacked layers, each parameterized by $\theta_l \in \Theta = \{\theta_1, \dots, \theta_L\}$. Let $\mathbf{x} \in \mathbb{R}^{N \times d}$ denote the input sequence (with N tokens and d -dimensional embeddings). The core layer function, \mathcal{T}_{θ_l} , typically consists of multi-head self-attention (MHSA) and feed-forward components, applied after layer normalization (LN):

$$H' = \text{MHSA}(\text{LN}(H)). \quad (12)$$

The forward computation for a standard Transformer model \mathcal{M}_{std} is expressed as a composition of the layer functions:

$$\mathcal{M}_{\text{std}}(\mathbf{x}) = \psi \circ \mathcal{T}_{\theta_L} \circ \dots \circ \mathcal{T}_{\theta_1} \circ \phi(\mathbf{x}), \quad (13)$$

where $\phi(\cdot)$ is an input embedding function and $\psi(\cdot)$ is an output head. Both the computational cost and the total number of parameters scale linearly with L .

C.2 Mixture-of-Experts

A Mixture-of-Experts layer replaces the standard dense feed-forward block with a sparse conditional computation module. For an input token h , a gating function routes the computation to a subset of K experts:

$$\text{ExpertOutput}(h) = \sum_{k=1}^K G_k(h) \cdot \mathcal{E}_{\theta_i^{(k)}}(h). \quad (14)$$

Load Balancing A critical challenge in training MoE models is expert collapse, where the router converges to selecting only a few experts. To mitigate this, a load balancing auxiliary loss \mathcal{L}_{aux} is typically added. For a batch of B tokens, let f_k be the fraction of tokens dispatched to expert k , and P_k be the average routing probability:

$$f_k = \frac{1}{B} \sum_{i=1}^B \mathbb{1}(\text{expert} = k), \quad P_k = \frac{1}{B} \sum_{i=1}^B G_k(h_i). \quad (15)$$

The auxiliary loss is defined as

$$\mathcal{L}_{\text{aux}} = N \sum_{k=1}^K f_k \cdot P_k, \quad (16)$$

which encourages uniform expert usage by maximizing the entropy of the routing distribution.

D Dynamics of Entropy-Adaptive Routing

In Top- p routing, the number of active experts N is determined dynamically. Let π denote the sorting permutation such that $p_{\pi(1)} \geq \dots \geq p_{\pi(K)}$. The active set size is formalized as:

$$N = \min \left\{ n \in \{1, \dots, K\} \mid \sum_{i=1}^n p_{\pi(i)} \geq p \right\} \quad (17)$$

This formulation creates an implicit functional dependency between the sparsity level N and the router entropy $\mathcal{H}(\mathbf{p})$. As the uncertainty varies between its extrema, N adapts:

$$N \approx \begin{cases} 1 & \text{as } \mathcal{H}(\mathbf{p}) \rightarrow 0 \quad (\text{High Confidence}) \\ \lceil p \cdot K \rceil & \text{as } \mathcal{H}(\mathbf{p}) \rightarrow \ln K \quad (\text{Max Entropy}) \end{cases} \quad (18)$$

Thus, the mechanism effectively ensembles more experts in high-entropy regimes while preserving sparsity in low-entropy scenarios.

E Stubborn Routing and Partition Errors

Similar to hallucinations in LLM sequence generation, MoE routers can become biased towards a specific expert subset despite poor fit. This is exacerbated by the accumulation of approximation errors in the router’s partition function Z_θ . Since Z_θ is computed only over the K existing experts, it fails to account for the “missing knowledge” that an ideal, infinite-capacity expert pool would provide. Consequently, a low-entropy routing state does not guarantee a reliable expert-token match; it merely indicates an absence of internal competition among the current experts, ignoring the possibility that the entire expert set may be ill-suited for the input.

F Extended Analysis of Routing States

State II: Functional Ambiguity. In this regime, the router possesses sufficient knowledge (low \mathcal{F}) but remains uncertain about a single vector (high \mathcal{H}). This indicates that the token represents a cross-disciplinary feature or “hub” concept that multiple experts are equally capable of handling. Unlike State IV, this uncertainty is aleatoric rather than epistemic, suggesting potential for collaborative activation or multi-routing strategies.

State IV: Cognitive Vacuum. This state represents the absence of learned structure. The model neither knows which expert to choose nor has any meaningful evidence for the input. Empirical observation suggests this state is dominant during the initial stages of training before expert specialization emerges, and persists in the final layers where residual stream variance often overwhelms the router’s feature extraction capabilities.

G Extended Analysis of Training Dynamics

In this section, we provide a detailed analysis of the training stability and thermodynamic evolution of the proposed methods, supplementing the main results.

G.1 Logit Magnitude and Evidence Potential

Logit sums ($\sum z_k$) serve as a proxy for *total potential evidence*. We observe a strong sensitivity to λ : high regularization ($\lambda = 10^{-2}$) suppresses magnitudes, acting as a temperature-like pressure preventing high-confidence commitment. Conversely,

weak regularization ($\lambda = 10^{-4}$) permits rapid evidence accumulation. Crucially, deeper layers exhibit consistent logit decay, signaling signal loss and a struggle to establish dominant expert presence in complex representation spaces.

G.2 Load Balancing as an Attachable Refinement

Our CLB is designed as an *attachable* module that can be layered on top of existing mechanisms—whether Auxloss or Auxfree—without modifying the underlying MoE architecture or increasing parameter count. In practice, CLB is activated mainly in the later stages of training. It selectively adjusts expert usage based on epistemic signals, refining the load distribution after the router and experts have already developed initial specialization.

The $\text{MaxViO}_{\text{batch}}$ curves over training steps (refer to Figure 8 in the main text and additional plots below) demonstrate that, once attached, CLB progressively reduces batch-wise imbalance. Crucially, it avoids the "early oscillation / late collapse" pattern frequently observed in conventional MoE training, where routers flip between unstable exploration and premature convergence to a few experts.

G.3 Thermodynamic Landscape and Routing Uncertainty

To investigate how UAR interacts with the router’s epistemic state, we visualize the evolution of routing uncertainty by tracking two key metrics across layers and training steps:

- **Helmholtz Free Energy (\mathcal{F}):** Acts as a proxy for absolute evidence.
- **Routing Entropy (\mathcal{H}):** Acts as a proxy for aleatoric uncertainty.

Visualizations indicate that enabling UAR qualitatively shifts the router’s thermodynamic trajectory: the energy "climb" becomes less pronounced, the high-level plateau is lowered, and fluctuations in both \mathcal{F} and \mathcal{H} are substantially reduced. These effects are consistent across two forms of intervention:

1. **Evidence-Triggered Expansion:** By modifying routing directly, the dynamic adjustment of active experts in high-energy regimes reduces the number of tokens trapped in "Cognitive Vacuum" and "Stubborn Bias" states.

Setting	Avg. ActE	FLOPs	Perplexity ↓
Top- k (Fixed k)	2.00	1.00x	23.93
Uncertainty = Entropy			
AE ($a = 0, b = 2$)	2.00	1.00x	23.93 (0.00)
AE ($a = 0.3, b = 2$)	2.15	1.00x	23.87 (-0.06)
AE ($a = 0.5, b = 2$)	2.45	1.01x	23.82 (-0.11)
Uncertainty = Energy			
AE ($a = 0.5, b = 2$)	2.55	1.01x	23.79 (-0.14)
AE ($a = 0.7, b = 2$)	2.80	1.02x	23.77 (-0.16)
AE ($a = 1.0, b = 2$)	3.10	1.03x	23.74 (-0.19)

Table 3: Ablation study of AE routing under different uncertainty measures (entropy/energy) and hyperparameter settings (a, b). We report average activated experts per token, FLOPs (normalized to Top- k), and validation perplexity. The values in green brackets show the improvement compared to the Top- k baseline.

This leads to systematically lower free energy and a smoother entropy decay.

2. **CLB Refinement:** When attached to Auxloss or Auxfree, CLB selectively applies regularization in high-uncertainty regions. This avoids unnecessary stochastic pressure on confident tokens, further stabilizing the router’s logit landscape.

Empirically, this joint reduction in routing energy and entropy variability indicates that UAR drives the MoE system toward a more calibrated uncertainty regime throughout training.

H Extended Experimental Analysis

In this section, we provide a granular breakdown of the thermodynamic evolution observed in Sec. 6.4, detailing the specific interactions between energy dynamics and routing behaviors.

H.1 The "Cognitive Vacuum" and Energy Suppression

A central premise of our framework is that standard Router Entropy is insufficient for detecting high epistemic uncertainty, particularly in deeper layers.

- **Evidence:** As illustrated in Fig. 8 (Left & Middle), the baseline model exhibits persistently high Free Energy in Layers 8–12, confirming the existence of a "Cognitive Vacuum" where the router lacks discriminative evidence.
- **Mechanism:** UAR transforms this dynamic. By utilizing Energy-AE, the model effectively

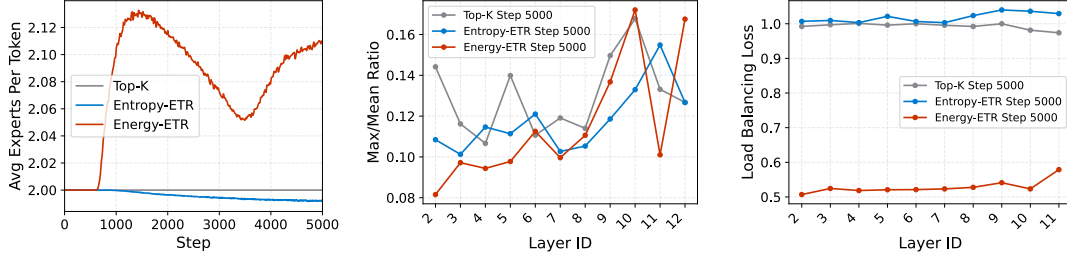


Figure 9: **Dynamic Expert Expansion and Load Balancing Efficiency.** **Left:** The average number of activated experts per token over time; Energy-AE dynamically expands the budget ($k_h > 2$) during the critical transition phase while Top- k remains fixed. **Middle:** The Max/Mean ratio of expert utilization across layers, where Energy-AE achieves the lowest ratio, indicating superior load uniformity. **Right:** Layer-wise auxiliary loss, validating that Energy-AE consistently minimizes the load balancing objective across all layers.

suppresses the energy floor in these deep layers. Instead of making low-confidence predictions, the router dynamically expands its capacity, "filling" the vacuum with computation via multi-expert ensembles.

H.2 Dynamic Expansion as a Phase Transition Buffer

Our theoretical analysis in Sec. 4 identified a critical "Energy Climb" during the Transition Stage (Phase B), characterized as *Frustrated Exploration*.

- **Adaptive Response:** Fig. 9 (Left) highlights a distinct peak in activated experts ($k_{eff} \approx 2.1$) specifically between steps 800 and 2000. This temporal window aligns perfectly with the identified "Symmetry Breaking" phase.
- **Interpretation:** The model automatically learns to expand k during periods of thermodynamic stress. This suggests that UAR enables the router to hedge bets during chaotic optimization landscapes. Crucially, as the model enters the Stabilization Phase (Phase C), k_{eff} naturally decays, validating the method's ability to navigate the exploration-exploitation trade-off autonomously.

H.3 Decoupling via Epistemic Masking

The "Diversity-Specialization Dilemma" suggests that uniform load balancing often impedes expert maturation.

- **Bias Evolution:** Fig. 7 (Right) reveals that the AuxFree baseline suffers from "Bias Explosion," where expert biases drift uncontrollably to satisfy crude routing heuristics.
- **CLB Efficacy:** The proposed CLB strategy (Blue/Grey lines) mitigates this by applying

Model Setting	L.2-Small	L.2-Middel	L.2-Large
<i>hidden size</i>	1024	1536	2048
<i>intermediate size</i>	2560	2560	4096
<i>attention heads</i>	32	32	32
<i>num kv heads</i>	32	16	32
<i>layers</i>	8	8	8
# Params	162M	230M	466M

Table 4: Detailed configuration, activation parameters, and total parameters of the models included in our study. L.2-162M represents the LLaMA-2 architecture model with 162M total parameters.

Epistemic Masking. Regularization is applied strictly to high-uncertainty tokens (State II/IV), preventing structural interference. This allows confident tokens (State I) to drive specialization without penalty, proving that uniformity can be improved without fighting the model's intrinsic drive for specialization.

1113
1114
1115
1116
1117
1118
1119