

Dynamic Large Concept Models: Latent Reasoning in an Adaptive Semantic Space

¹ByteDance Seed, ²University of Manchester, ³Mila - Quebec AI Institute, ⁴Tsinghua University, ⁵M-A-P

Abstract

Large Language Models (LLMs) apply uniform computation to all tokens, despite language exhibiting highly non-uniform information density. This token-uniform regime wastes capacity on locally predictable spans while under-allocating computation to semantically critical transitions. We propose **Dynamic Large Concept Models (DLCM)**, a hierarchical language modeling framework that learns semantic boundaries from latent representations and shifts computation from tokens to a compressed concept space where reasoning is more efficient. DLCM discovers variable-length concepts end-to-end without relying on predefined linguistic units. Hierarchical compression fundamentally changes scaling behavior. We introduce the first **compression-aware scaling law**, which disentangles token-level capacity, concept-level reasoning capacity, and compression ratio, enabling principled compute allocation under fixed FLOPs. To stably train this heterogeneous architecture, we further develop a **decoupled μP parametrization** that supports zero-shot hyperparameter transfer across widths and compression regimes. At a practical setting ($R = 4$, corresponding to an average of four tokens per concept), DLCM reallocates roughly one-third of inference compute into a higher-capacity reasoning backbone, achieving a **+2.69% average improvement** across 12 zero-shot benchmarks under matched inference FLOPs.

Date: January 1, 2026

Correspondence: xingwei.qu@bytedance.com, zhangge.eli@bytedance.com

1 Introduction

Large Language Models (LLMs) have achieved remarkable success across natural language understanding, reasoning, and generation tasks. Despite differences in scale and training data, nearly all state-of-the-art models share a common architectural assumption: language is processed uniformly at the token level, with identical depth and computation applied to every position in the sequence.

This assumption stands in sharp contrast to the structure of natural language. Information density is highly non-uniform: long spans of locally predictable tokens are interspersed with sparse but semantically critical transitions where new concepts are introduced and reasoning difficulty concentrates. Yet standard LLMs expend full computation on both regimes alike, resulting in substantial redundancy and systematic misallocation of model capacity.

More fundamentally, this inefficiency reflects a limitation of token-level modeling itself. Reasoning is inherently hierarchical: humans reason over abstract units such as ideas or concepts before committing to surface realizations. Token-level autoregressive models, however, lack any explicit abstraction mechanism and are forced to repeatedly infer high-level structure implicitly at every layer, solely through next-token prediction.

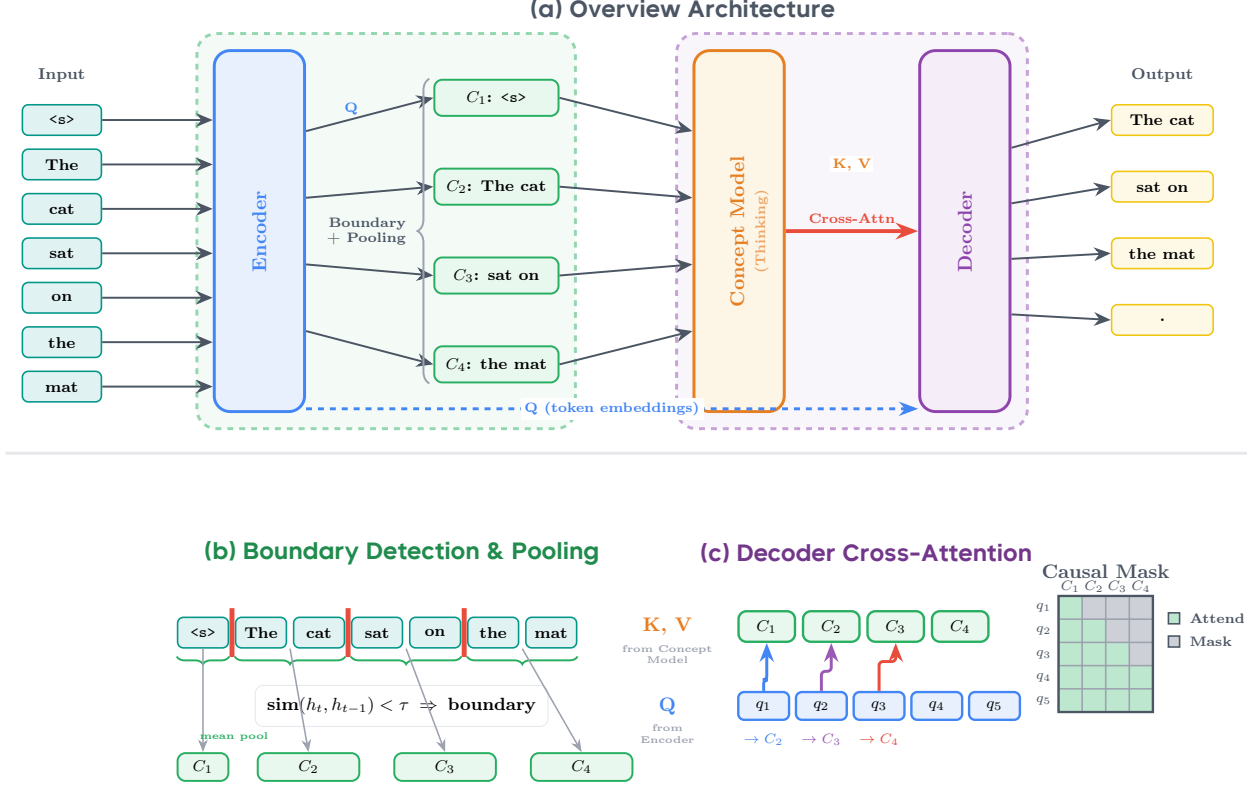


Figure 1 Overview Structure of DLCM.

Prior work has explored relaxing this constraint, but with important limitations. Latent reasoning approaches perform inference in continuous hidden spaces without explicit token generation, while sentence-level concept models rely on fixed, human-defined segmentation. Neither enables models to learn where semantic computation should be concentrated.

We argue that effective reasoning requires a learned intermediate granularity: neither raw tokens nor predefined sentences, but variable-length semantic concepts discovered directly from representation space. Based on this insight, we propose **Dynamic Large Concept Models (DLCM)**, a hierarchical next-token prediction framework that dynamically segments token sequences into concepts, performs deep reasoning in a compressed concept space, and reconstructs token-level predictions through causal cross-attention.

This design separates what to reason about from how to reason. By learning semantic boundaries end-to-end and relocating computation from redundant token processing to concept-level reasoning, DLCM enables adaptive compute allocation aligned with information density.

At the other extreme, H-NET [11] demonstrates the promise of learned boundary detection with adaptive compute allocation, but operates at the byte level and has not been validated against standard Next Token Prediction (NTP) baselines in modern LLM pipelines.

Our work bridges these gaps by introducing **latent reasoning at the concept level**. Throughout this paper, we use the term concept to denote variable-length latent segments discovered from representation space, rather than linguistically predefined semantic units. The key insight is that effective reasoning requires neither token-level granularity (too fine, computationally wasteful) nor sentence-level granularity (too coarse, inflexible), but rather semantically coherent concepts whose boundaries are learned end-to-end from data. We propose **DLCM**, a hierarchical architecture that implements this insight through a four-stage pipeline (Figure 1):

1. **Encoding:** A lightweight encoder processes raw tokens to extract fine-grained representations.

2. **Dynamic Segmentation:** A learned boundary detector identifies semantic breakpoints by measuring local dissimilarity between adjacent token representations. Unlike LCM’s fixed sentence boundaries, these boundaries emerge from the model’s own latent space through end-to-end optimization.
3. **Concept-Level Reasoning:** Tokens within each segment are pooled into unified concept representations. A high-capacity transformer then performs deep reasoning exclusively on this compressed concept sequence—where the majority of computation occurs.
4. **Token-Level Decoding:** A decoder reconstructs token-level predictions by attending to the reasoned concepts via a causal cross-attention mechanism.

This design explicitly decouples what to think about (concept formation via learned boundaries) from how to think (reasoning in compressed latent space), enabling the model to allocate computation adaptively based on semantic structure rather than surface token count. As our results demonstrate, this structural bias makes the model exceptionally proficient at handling high-information, low-predictability tokens that mark the beginning of new concepts.

Our main contributions are as follows:

- **Concept-level latent reasoning with learned boundaries.** We propose DLCM, a hierarchical next-token prediction architecture that discovers variable-length semantic concepts from latent representations and performs deep computation in a compressed concept space.
- **Compression-aware scaling law for hierarchical LMs.** We derive a scaling law $L(N, D, R, P)$ that explicitly models the interaction between total parameters N , data D , compression ratio R , and concept-backbone allocation P , enabling principled selection of architecture under equal-FLOPs constraints. [9]
- **Decoupled μP for heterogeneous modules (why different LR / init variance).** We demonstrate that the Maximal Update Parametrization (μP) can be effectively adapted to our heterogeneous architecture to prevent training instability and ensure optimal performance at scale. Specifically, we identify that due to the decoupled widths of our model, the learning rates for the token-level components, concept backbone, and embeddings must be adjusted independently. Empirically, we confirm that the optimal effective learning rate for each component scales inversely with its specific width ($\eta \propto \text{width}^{-1}$), a finding that aligns with theoretical predictions for uniform models but is verified here in a non-uniform setting.
- **Compute redistribution yields reasoning gains at lower FLOPs.** With $R = 4$, DLCM reduces FLOPs by up to 34% while reallocating capacity into a larger reasoning backbone, improving average accuracy by 2.69% on 12 zero-shot benchmarks, with the largest gains on reasoning-dominant tasks.

2 Related Work

2.1 From Latent Reasoning to Concept-Level Language Modeling

Recent research has explored performing reasoning at higher levels of abstraction than individual tokens, offering both computational efficiency and new modeling capabilities. **Latent reasoning** frameworks perform reasoning entirely within continuous hidden state spaces rather than through explicit token generation [6]. In the COCONUT framework, the model’s hidden state from one reasoning step feeds directly into subsequent steps without generating intermediate tokens [6]. This approach offers significant computational advantages over methods like Chain-of-Thought prompting, which require generating hundreds of intermediate tokens. More importantly, continuous representations can encode multiple potential reasoning paths in superposition, enabling parallel exploration of the solution space [6]. However, this comes at the cost of interpretability and can struggle with tasks requiring precise symbolic manipulation.

Building upon latent reasoning principles, the **Large Concept Model (LCM)** framework operates at an intermediate level: reasoning on sentence-level "concepts" rather than tokens [17]. LCMs use a three-stage pipeline: (1) a frozen encoder maps sentences to fixed-size embeddings in a semantic space (e.g., SONAR, supporting 200 languages), (2) a transformer performs autoregressive prediction in this concept space using

diffusion or quantization techniques adapted from computer vision, and (3) a frozen decoder reconstructs token-level output [17]. This approach combines key advantages: like latent reasoning, it operates in continuous space with substantial efficiency gains (10× sequence length reduction), but unlike pure latent reasoning, each concept remains interpretable as a decodable sentence [17]. Remarkably, LCMs demonstrate zero-shot multilingual transfer—models trained only on English can generate in 200+ languages by leveraging the language-agnostic semantic space [17]. However, the LCM framework faces significant limitations. First, it requires pre-training separate encoder and decoder models on massive multilingual data before the LCM itself can be trained, creating a scalability bottleneck [17]. Second, and more fundamentally, the sentence-level granularity is a fixed human prior—the model must accept predetermined sentence boundaries rather than learning task-optimal segmentation. This rigidity prevents the model from adapting its conceptual granularity to different domains or tasks. Our DCLM architecture addresses both limitations through end-to-end training with dynamic, learnable boundary detection that discovers optimal chunking strategies directly from the data.

2.2 Dynamic Compute Allocation in Language Models

Standard LLMs allocate uniform computation to every token, ignoring that some tokens (e.g., predictable function words) require minimal processing while others (e.g., concept boundaries) demand more effort [8, 13]. Recent work explores adaptive allocation mechanisms. The Universal Transformer [3] introduced recurrence in depth, applying the same transformation block repeatedly with a learned halting mechanism that determines when each position has been sufficiently refined. Mixture of Experts (MoE) models [12, 15] achieve conditional computation by routing each token to a subset of expert sub-networks. However, these approaches focus on parameter efficiency and scaling rather than fundamentally addressing the information density problem.

H-NET [11] directly addresses adaptive allocation through learned boundary detection. The model predicts semantic boundaries by analyzing local patterns (e.g., similarity between consecutive hidden states), segments the sequence into variable-length chunks, and processes the compressed chunk representations hierarchically. Critically, boundary detection is differentiable and trained end-to-end, allowing task-appropriate chunking strategies to emerge [11]. This yields substantial gains: learned boundaries align with linguistic structures even without supervision, and compression (4-8× reduction) translates to quadratic attention savings [11]. The approach implicitly allocates more computation to high-information boundaries where new concepts begin—precisely where prediction is most difficult. However, H-NET’s primary focus is on efficient representation through hierarchical bit-level modeling rather than token-level generation in state-of-the-art autoregressive LLMs [11]. This leaves unaddressed the critical problem of computational waste in modern decoder-only language models, where every token—regardless of its predictability or information content—receives identical processing through the full model depth. Our DCLM architecture bridges this gap by adapting H-NET’s dynamic boundary detection principles to the token-level generation paradigm of current LLMs [1, 18]. By segmenting sequences into concept chunks and performing compressed reasoning before token-level decoding, DCLM enables adaptive computation allocation in the exact architectural context where it matters most: next-token prediction in large-scale autoregressive models.

3 Methodology

We now describe the technical details of DCLM. The overall architecture is illustrated in Figure 1.

3.1 Overview

DCLM processes a token sequence through four stages: (1) **Encoding** extracts fine-grained token representations; (2) **Dynamic Segmentation** identifies semantic boundaries and pools tokens into concepts; (3) **Concept-Level Reasoning** performs deep computation on the compressed sequence; and (4) **Token-Level Decoding** reconstructs predictions by attending to reasoned concepts. We formalize this as:

$$\mathbf{H} = \mathcal{E}(\mathbf{x}) \quad (\text{Encoding}) \quad (1)$$

$$\mathbf{C} = \Phi(\mathbf{H}) \quad (\text{Segmentation \& Pooling}) \quad (2)$$

$$\mathbf{Z} = \mathcal{M}(\mathbf{C}) \quad (\text{Concept Reasoning}) \quad (3)$$

$$\hat{\mathbf{y}} = \mathcal{D}(\Psi(\mathbf{H}, \mathbf{Z})) \quad (\text{Decoding}) \quad (4)$$

where \mathcal{E} is the encoder, Φ is the segmentation-pooling operation, \mathcal{M} is the concept-level transformer, \mathcal{D} is the decoder, and Ψ is the cross-attention expansion defined in Eq. 14.

3.2 Encoding

The encoder \mathcal{E} is a standard causal Transformer that processes raw tokens $\mathbf{x} = [x_1, \dots, x_L]$ to produce fine-grained representations $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_L] \in \mathbb{R}^{L \times d_{\text{token}}}$. These representations capture local contextual information and serve as the basis for both boundary detection and final token-level decoding.

3.3 Dynamic Segmentation

While boundary scores are learned end-to-end from latent representations, we intentionally decouple the discrete segmentation decision from the language modeling loss to avoid optimization interference. This design trades full end-to-end discreteness for training stability and controllable compression, which we find essential at scale.

3.3.1 Boundary Detection

Our key hypothesis is that transitions between distinct concepts are marked by significant shifts in the latent feature space. We detect these “semantic breaks” by measuring local dissimilarity between adjacent tokens.

Given encoder outputs \mathbf{H} , we project each token into a query-key space of dimension d_{scan} :

$$\mathbf{q}_t = \mathbf{W}_q \mathbf{h}_t, \quad \mathbf{k}_t = \mathbf{W}_k \mathbf{h}_t \quad (5)$$

The boundary probability $p_t \in [0, 1]$ is computed as the normalized dissimilarity:

$$p_t = \frac{1 - \cos(\mathbf{q}_{t-1}, \mathbf{k}_t)}{2} = \frac{1}{2} \left(1 - \frac{\mathbf{q}_{t-1}^\top \mathbf{k}_t}{\|\mathbf{q}_{t-1}\|_2 \|\mathbf{k}_t\|_2} \right) \quad (6)$$

We enforce $p_1 = 1$ so that the first token always starts a new concept.

Discrete Sampling. While p_t is continuous, downstream processing requires discrete segment assignments. We adopt different strategies for training and inference:

- **Training:** We sharpen probabilities by temperature α , then sample $b_t \sim \text{Bernoulli}(p_t^{\text{sharp}})$ to encourage exploration.
- **Inference:** We use a **hard thresholding** rule: $b_t = \lfloor p_t \geq 0.5 \rfloor$.

3.3.2 Concept Formation via Pooling

Given boundary indicators $\mathbf{b} = [b_1, \dots, b_L]$, we partition tokens into M contiguous segments S_1, \dots, S_M . Each segment is compressed into a single concept representation via mean pooling, followed by a projection to the concept dimension d_{concept} :

$$\mathbf{c}_k^{\text{raw}} = \frac{1}{|S_k|} \sum_{t \in S_k} \mathbf{h}_t, \quad \mathbf{c}_k = \mathbf{W}_{\text{up}} \mathbf{c}_k^{\text{raw}} \quad (7)$$

where $\mathbf{W}_{\text{up}} \in \mathbb{R}^{d_{\text{concept}} \times d_{\text{token}}}$ aligns the feature space. The resulting concept sequence $\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_M]$ has length $M \ll L$.

3.3.3 Adaptive Compression via Global Load Balancing

Similar to H-Net [11], natural language exhibits varying information density. To enable content-adaptive compression while maintaining a target ratio R (e.g., $R = 4$ means 4 tokens per concept on average), we impose constraints at the global batch level rather than per-sequence.

Let \mathcal{T} denote all tokens across the distributed batch. We track:

$$G_{\text{global}} = \frac{1}{|\mathcal{T}|} \sum_{(i,t) \in \mathcal{T}} p_{i,t} \quad (\text{expected boundary rate}) \quad (8)$$

$$F_{\text{global}} = \frac{1}{|\mathcal{T}|} \sum_{(i,t) \in \mathcal{T}} b_{i,t} \quad (\text{actual boundary rate}) \quad (9)$$

These statistics are synchronized across ranks via **AllReduce**. We optimize an auxiliary loss:

$$\mathcal{L}_{\text{aux}} = \frac{R}{R-1} [(R-1) \cdot F_{\text{global}} \cdot G_{\text{global}} + (1 - F_{\text{global}}) \cdot (1 - G_{\text{global}})] - 1 \quad (10)$$

This encourages the global compression rate to converge to $1/R$ while allowing local fluctuation. We refer to this globally regularized segmentation mechanism as the *Global Parser*, which serves as a critical component for enabling content-adaptive granularity.

3.4 Concept-Level Reasoning

The concept-level transformer \mathcal{M} is the computational core of DLCM. Operating on the compressed sequence $\mathbf{C} \in \mathbb{R}^{M \times d_{\text{concept}}}$, it performs deep reasoning with significantly reduced attention complexity.

\mathcal{M} is a standard causal Transformer with L_{concept} layers. Because concepts represent semantic units rather than surface tokens, this module can focus on high-level reasoning without being distracted by low-level token prediction. The output $\mathbf{Z} = \mathcal{M}(\mathbf{C})$ contains enriched concept representations.

3.5 Token-Level Decoding

The decoder reconstructs token-level predictions by attending to the reasoned concepts. This involves two components: concept smoothing and causal cross-attention.

3.5.1 Concept Smoothing

Hard pooling can introduce discretization artifacts at segment boundaries. We apply a lightweight smoothing module \mathcal{S} to integrate adjacent concepts:

$$\tilde{\mathbf{Z}} = \mathcal{S}(\mathbf{Z}) \quad (11)$$

3.5.2 Causal Cross-Attention

The decoder \mathcal{D} generates token predictions by querying the smoothed concepts. Crucially, we enforce causality so token t can only attend to concepts formed up to index $j(t) = \sum_{i=1}^t b_i$.

To handle the heterogeneous architecture ($d_{\text{token}} \neq d_{\text{concept}}$), the cross-attention mechanism projects queries from the encoder space and keys/values from the concept space into a common head dimension d_{head} :

$$\mathbf{Q} = \mathbf{H}\mathbf{W}_Q, \quad \text{where } \mathbf{W}_Q \in \mathbb{R}^{d_{\text{token}} \times d_{\text{head}}} \quad (12)$$

$$\mathbf{K} = \tilde{\mathbf{Z}}\mathbf{W}_K, \quad \mathbf{V} = \tilde{\mathbf{Z}}\mathbf{W}_V, \quad \text{where } \mathbf{W}_{K,V} \in \mathbb{R}^{d_{\text{concept}} \times d_{\text{head}}} \quad (13)$$

The attention output is computed as:

$$\Psi(\mathbf{H}, \mathbf{Z}) = \text{Softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_{\text{head}}}} + \mathbf{M} \right) \mathbf{V}\mathbf{W}_O + \mathbf{H} \quad (14)$$

where $\mathbf{W}_O \in \mathbb{R}^{d_{\text{head}} \times d_{\text{token}}}$ projects the result back to the token dimension for the residual connection.

3.6 Training Objective

The total loss combines next-token prediction with adaptive compression:

$$\mathcal{L} = \mathcal{L}_{\text{CE}} + \lambda \mathcal{L}_{\text{aux}} \quad (15)$$

where \mathcal{L}_{CE} is cross-entropy on output tokens and \mathcal{L}_{aux} is the load-balancing loss.

4 Implementation Details

Packed Sequence Training. We adopt the Variable Length (VarLen) approach from FlashAttention [2] to ensure global compression statistics are computed over diverse tokens.

QK Normalization. Bridging token-level and concept-level representations with different statistical properties can cause training instability. Following [4, 7], we apply RMSNorm to queries and keys before attention:

$$\mathbf{Q}' = \text{RMSNorm}(\mathbf{Q}), \quad \mathbf{K}' = \text{RMSNorm}(\mathbf{K}) \quad (16)$$

4.1 Efficient Cross-Attention via Concept Replication

The decoder’s cross-attention mechanism (Section 3.5) presents a significant implementation challenge. Mathematically, tokens must attend to concepts with variable-length mappings ($L \times M$), creating irregular attention patterns. As illustrated in Figure 2, when tokens $\{t_1\}$ belong to concept c_1 , and tokens $\{t_2, t_3\}$ belong to c_2 , the resulting attention mask effectively has a "ragged" boundary.

Implementing this directly with **Flex Attention** incurs significant overhead from dynamic mask generation and irregular memory access patterns. To address this, we adopt a **concept replication strategy** to bridge the gap between the theoretical $L \times M$ formulation and hardware-friendly $L \times L$ kernels. Analogous to Grouped Query Attention (GQA), we expand concepts via `repeat_interleave` to match token positions.

Specifically, for each token t_i belonging to concept c_j , we replicate the concept feature c_j at position i in the key/value sequence:

$$\tilde{\mathbf{K}} = \text{repeat_interleave}(\mathbf{K}, \text{segment_lengths}), \quad \tilde{\mathbf{V}} = \text{repeat_interleave}(\mathbf{V}, \text{segment_lengths}) \quad (17)$$

This transformation aligns the Key/Value length with the Query length (L), enabling the use of **Flash Attention with Variable Length (Varlen)**. This allows us to leverage highly optimized CUDA kernels designed for standard causal masking, treating the problem as a specialized form of self-attention where keys and values are locally constant within each concept segment.

4.2 Performance Benchmarks

We benchmark the efficiency of our concept replication strategy against Flex Attention using independent kernel profiling. Note that the hidden sizes listed here (1024, 2048, 4096) are standard benchmarking dimensions and do not necessarily match the specific architectural dimensions ($d_{\text{token}}, d_{\text{concept}}$) of DLCM, as the primary goal is to demonstrate algorithmic scalability.

Table 6 provides detailed performance. To more intuitively analyze the performance trends, we have plotted the speedup ($T_{\text{flex}}/T_{\text{flash}}$) in Figure 9.

4.3 Key Observations and Analysis

We can draw three key conclusions from these results:

- **Consistent Performance Advantage:** In all tested configurations, Flash Attention Varlen with concept replication significantly outperforms Flex Attention. The speedup ranges from **1.26×** to **1.73×**, validating the efficacy of the “memory-for-computation” trade-off.

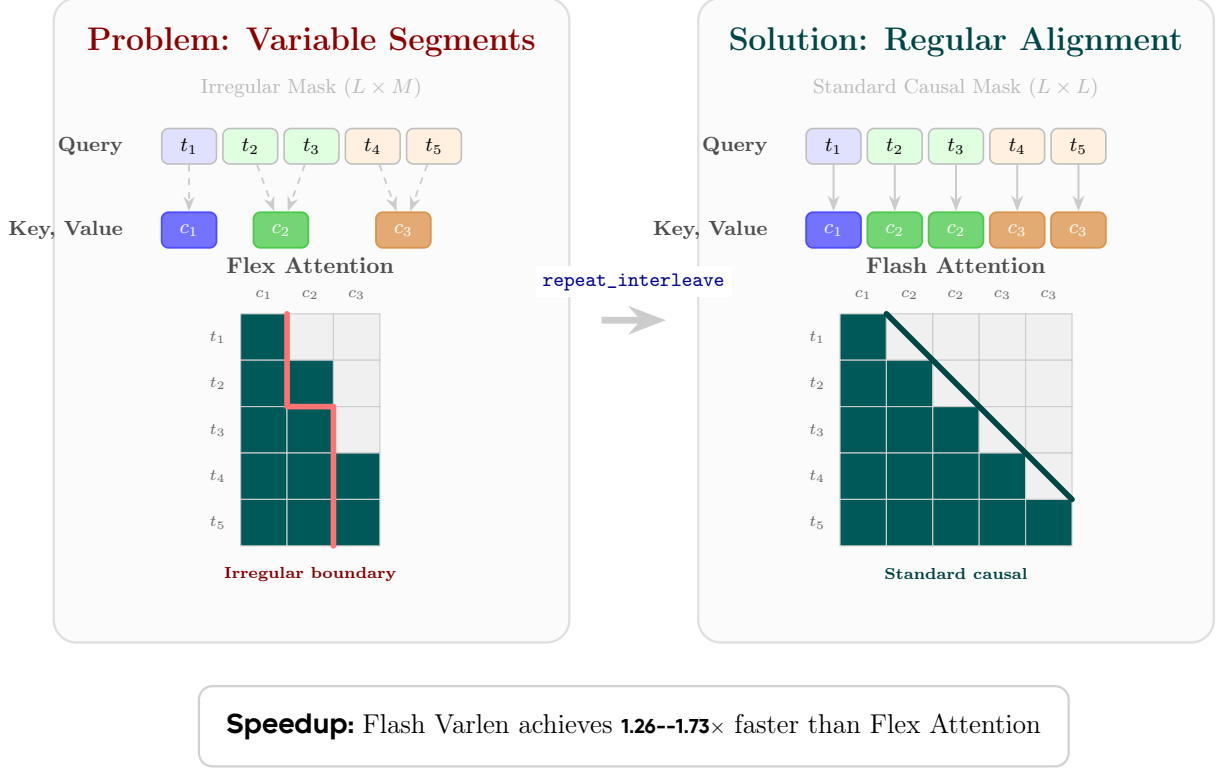


Figure 2 Cross-Attention Optimization via Concept Replication. *Left:* The decoder’s cross-attention creates an irregular $L \times M$ mask due to variable token-to-concept mappings. *Right:* By replicating concepts via `repeat_interleave` to match token positions, we obtain a standard $L \times L$ causal mask, enabling optimized Flash Attention kernels.

- **Insensitivity to Hidden Size:** As shown in Figure 9, the three lines representing different hidden sizes (1024, 2048, and 4096) are nearly identical. This indicates that the performance bottleneck is dominated by the memory access patterns of the attention mechanism, not the computational complexity of the hidden dimension. Flash Varlen’s optimized, regular memory access pattern remains stable across various model widths.
- **Superior Scalability with Sequence Length:** The most critical finding is that Flash Varlen’s performance advantage **scales with increasing sequence length**. At a 2K sequence length, the average speedup is $\sim 1.44\times$. When the sequence length increases 8-fold to 16K, the average speedup climbs to $\sim 1.70\times$, peaking at **1.73 \times** for a hidden size of 2048.

This scalability trend strongly suggests that the overhead from Flex Attention’s dynamic mask generation and irregular memory access patterns grows faster than the computational cost. Conversely, despite its increased memory footprint for the K/V cache, Flash Varlen’s highly optimized kernel and regular causal access pattern prove far more efficient, especially at longer sequences.

5 Data

To ensure experimental reproducibility, we build our corpus entirely from open-source data and tokenize it using the DeepSeek-v3 [14] tokenizer. Our corpus spans multiple domains, including web text (English and Chinese), mathematics, and code, forming a comprehensive foundation for core language understanding across linguistic, factual, and reasoning abilities.

The composition of our pretraining data is intentionally designed to serve two critical objectives. First, it balances breadth and specialization: web text provides broad natural language coverage, while mathematics

and code enhance structured reasoning. Second, and more importantly for our architecture, this diversity is essential for learning robust dynamic segmentation. By exposing the model to domains with drastically different information densities (e.g., highly structured code syntax vs. verbose natural language prose), we force the learned boundary predictor to discover content-adaptive segmentation strategies that generalize across diverse tasks. English and Chinese web text are weighted more heavily to ensure multilingual alignment, while specialized datasets like MegaMath-Web and OpenCoder-Pretrain are included to fine-tune the model’s handling of high-entropy transitions.

To demonstrate the architectural benefits of DLCM rather than gains from data curation, we do not apply aggressive filtering; instead, we use data whose quality aligns with standard open-source corpora. Table 1 summarizes the statistics.

Table 1 Statistics of the pretraining data.

| Data Source | Ratio | Tokens (B) |
|--------------------------------|-------------|--------------|
| Nemotron-CC [16] (English Web) | 50% | 500 |
| MAP-CC [5] (Chinese Web) | 25% | 250 |
| OpenCoder-Pretrain [10] | 15% | 150 |
| MegaMath-Web [20] | 10% | 100 |
| Total | 100% | 1,000 |

6 Scaling Laws for DLCM

To determine the optimal architecture and hyperparameters for DLCM, we conduct a comprehensive exploration using scaling laws. We first introduce a decoupled optimization strategy to handle the heterogeneous nature of our architecture, followed by the mathematical formulation of our scaling objectives.

6.1 Decoupled μ P for Heterogeneous Architectures

6.1.1 Formulation of μ P

To ensure consistent feature learning dynamics across varying scales and compression rates, we adopt the Maximal Update Parametrization (μ P). Unlike standard transformers with uniform width, our architecture requires decoupled scaling for the token-level components (\mathcal{E}, \mathcal{D}) with width d_{token} and the concept-level backbone (\mathcal{M}) with width d_{concept} .

We define distinct width multipliers relative to a base width d_{base} :

$$s_{\text{token}} = \frac{d_{\text{token}}}{d_{\text{base}}}, \quad s_{\text{concept}} = \frac{d_{\text{concept}}}{d_{\text{base}}} \quad (18)$$

Following standard μ P practice, we adjust initialization variances and optimization hyperparameters separately for each component group:

- **Initialization:** All hidden linear weights $W \in \mathbb{R}^{d_{\text{out}} \times d_{\text{in}}}$ are initialized with variance $\sigma_{\text{base}}^2 \cdot s^{-1}$, where $s \in \{s_{\text{token}}, s_{\text{concept}}\}$ corresponds to the layer’s width. Embedding weights use fixed σ_{base}^2 .
- **Learning Rates:** To stabilize feature learning, hidden layer learning rates are scaled inversely to width:

$$\eta_{\mathcal{E}, \mathcal{D}} = \eta_{\text{token}}^{\text{base}} \cdot s_{\text{token}}^{-1} \quad (19)$$

$$\eta_{\mathcal{M}} = \eta_{\text{concept}}^{\text{base}} \cdot s_{\text{concept}}^{-1} \quad (20)$$

Biases and embedding weights retain the fixed learning rate $\eta_{\text{others}}^{\text{base}}$.

- **Output Scaling:** To ensure the logits remain $O(1)$ effectively, the final decoder projection W_{unemb} is scaled during the forward pass:

$$\text{logits} = \frac{1}{s_{\text{token}}} \cdot (\mathbf{h}_{\text{final}} W_{\text{unemb}}^{\top}) \quad (21)$$

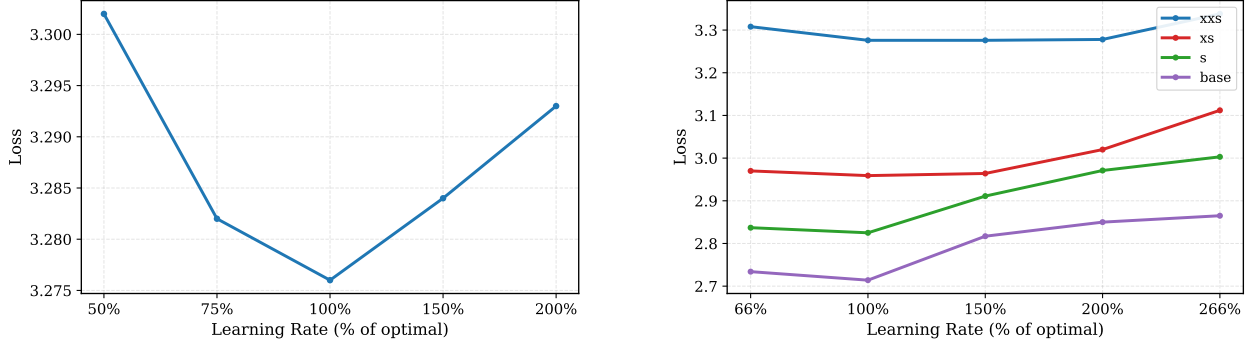


Figure 3 Hyperparameter tuning and transfer under μP . **Left:** We sweep $\eta_{\text{concept}}^{\text{base}}$ while fixing $\eta_{\text{token}}^{\text{base}}$ at its μP -predicted value on the 87M proxy model. The loss curve shows a well-defined minimum near 100%. **Right:** We jointly scale both base learning rates by the same factor and observe consistent minima across model sizes (87M–834M), validating the zero-shot transferability of μP -derived hyperparameters.

- **Optimizer Stability:** The AdamW ϵ parameter for each layer is scaled by s^{-1} , matching the respective component width.

6.1.2 Hyperparameter Tuning and Verification

Following the protocol proposed by Yang et al. [19], we adopted a two-stage strategy: tuning hyperparameters on a small proxy model and verifying their transferability on larger scales.

Proxy Model Tuning. We performed coordinate descent on the base learning rates using a proxy model with 87M parameters. For each hyperparameter group, we iteratively swept over a multiplicative grid of $\{0.5, 0.75, 1.5, 2.0\}$ relative to the current best value until the validation loss stabilized. Empirically, we observed that the optimal base learning rates for the token and concept components were approximately equal ($\eta_{\text{concept}}^{\text{base}} \approx \eta_{\text{token}}^{\text{base}}$). We show the result of tuning $\eta_{\text{concept}}^{\text{base}}$ while keeping $\eta_{\text{token}}^{\text{base}}$ in Figure 3, which means the actual learning rates depend on the ratio of the widths between the token and concept components in this unequal-width model. This consistency suggests that the explicit width-dependent scaling factors defined previously successfully account for the structural differences between components, stabilizing the effective learning rates across different widths.

Transfer Verification. To validate the zero-shot hyperparameter transfer, we trained larger models (274M, 468M, 834M parameters) using the optimal η^{base} values derived from the 87M proxy. Then we perturbed the predicted learning rates by simultaneously scaling $\eta_{\text{token}}^{\text{base}}$ and $\eta_{\text{concept}}^{\text{base}}$. As shown in Figure 3, deviating from the μP -predicted learning rates resulted in degraded performance, confirming that the optimal hyperparameters found on the proxy model transfer effectively to larger scales without further tuning. Overall, this confirms that μP effectively stabilizes training for our unequal-width architecture, provided that the learning rates for the token and concept components are decoupled and scaled inversely to their respective widths. This finding extends standard scaling laws to heterogeneous designs, ensuring consistent optimality across scales.

6.2 Scaling Law Formulation

6.2.1 Experimental Setup

To validate our scaling hypotheses, we constructed a grid of models by varying the concept-layer parameter ratio $P \in \{30\%, 50\%, 70\%\}$ and the compression ratio $R \in \{2, 4, 8\}$. Models were trained on a budget of 200B tokens, resulting in three primary model scales for analysis: **Small** (274M), **Medium** (468M), and **Large** (833M). All scaling exponents ($\delta_1, \delta_2, \gamma$) are shared globally across all model scales and compression ratios, and are fitted once using the joint training trajectories. Only scale-independent offset terms are allowed to

vary across configurations. This design constrains the degrees of freedom of the model and avoids post-hoc overfitting to individual scales.

6.2.2 Mathematical Formulation

We extend the Chinchilla scaling framework [8] by introducing (i) compression-aware behaviour and (ii) architectural decomposition. The resulting loss law is:

$$L(N, D, R, P) = E_0 + \frac{A_{\text{token}}}{(N(1 - P) + t_{\text{token}})^{\delta_1}} + \frac{A_{\text{concept}} R^\gamma}{(NP + t_{\text{concept}})^{\delta_2}} + \frac{A_{\text{data}}}{(D + t_{\text{data}})^\alpha}. \quad (22)$$

where N is total parameters, D is dataset size, R is compression ratio, P is concept-layer parameter ratio, and E_0 is the irreducible loss floor. This formulation disentangles token-processing efficiency, concept-processing efficiency (controlled by exponent γ), and data scaling.

6.2.3 Decay-Phase Power Law

Since our training protocol involves Weight-Sharing-and-Decay (WSD), we explicitly model the late-stage regime. We fit a simplified decay law to the fractional loss reduction Δ_{decay} in the 90%–99% token window:

$$\Delta_{\text{decay}} = k L_{\text{stable}}^a R^b N^c \quad (23)$$

Obtained via log-linear regression, this model achieves $R^2 = 0.93$, accurately predicting late-stage loss drops across all scales.

As shown in Figure 4 and Figure 5, our methodology—incorporating a tail-focused sampling strategy and weighting late-token regions—ensures the law generalizes reliably across both architectural and data scales.

6.3 Optimal Configuration Analysis

6.3.1 Architectural Efficiency

Figure 6(a) illustrates the Loss/FLOPs efficiency across different backbone proportions P and compression ratios R .

Selection of Compression Ratio ($R = 4$): While higher compression rates offer theoretical FLOPs savings, we empirically selected $R = 4$ as the primary configuration. This decision is driven by the granularity of concept compression; as discussed in Appendix A, a compression ratio of 4 aligns better with the intuitive semantic segmentation of tokens into concepts, offering the best balance between training stability and computational efficiency.

6.4 Scaling Law Validation and Verification

We developed a unified scaling-law estimator by jointly modeling the full-training loss trajectory and the late-stage decay behavior. To ensure robustness, we adopted a tail-focused sampling strategy that emphasizes curvature near convergence, performing fits on 100B-token trajectories while validating against 1T-token limits. As shown in Figure 4 and Figure 5, our estimator maintains a fitting error below 0.05 across the entire window.

This high-fidelity fitting allows for a critical verification of our architectural properties: our scaling law yields an effective compute multiplier prediction of approximately **1.4**. This value aligns closely with the standard baseline factor of **1.34**. This consistency confirms that our theoretical projections are grounded in established empirical norms and that the architecture scales predictably under the proposed law.

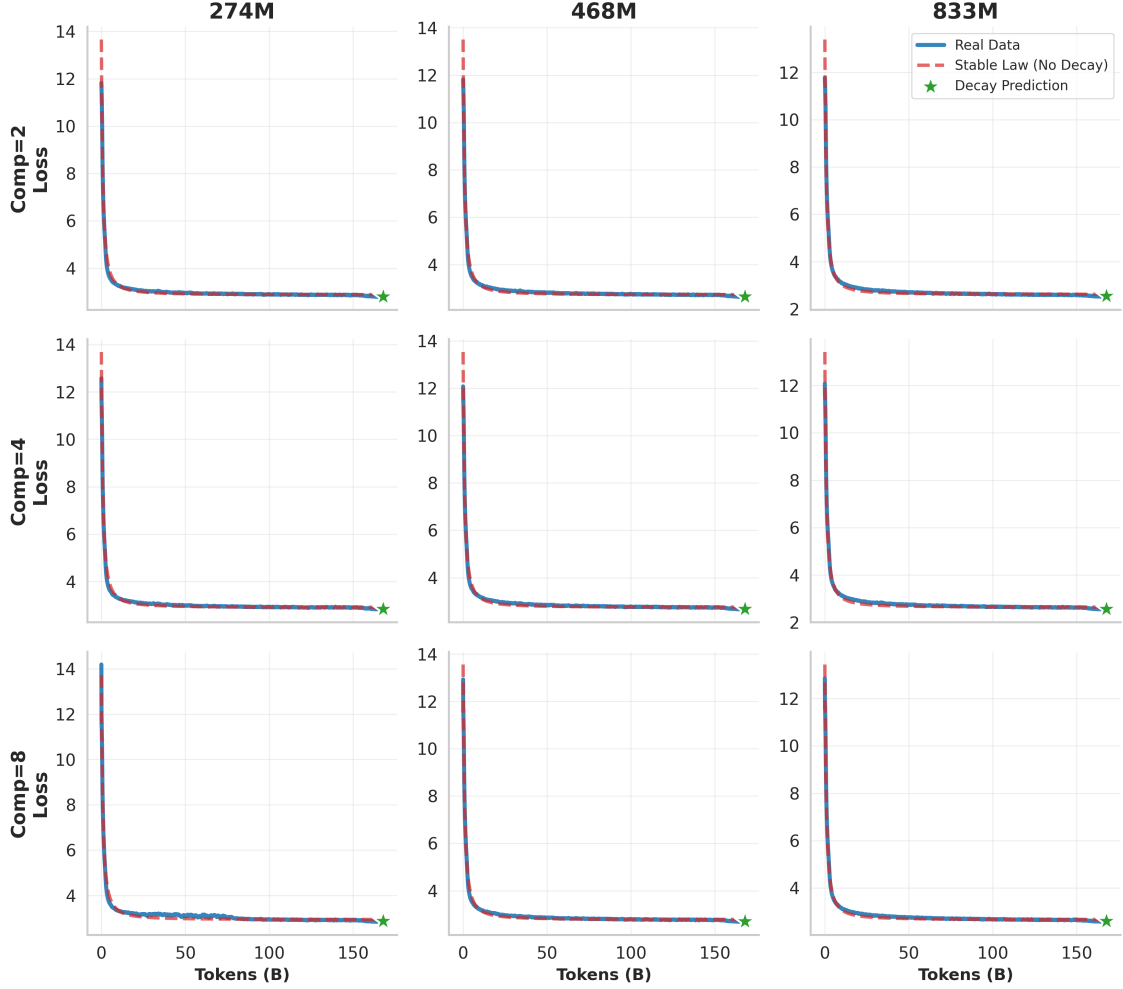


Figure 4 Full training trajectory fit. Comparison between predicted loss (Equation 22) and empirical loss across model sizes (274M–833M), compression factors $R \in \{2, 4, 8\}$, and training budgets. The joint fit achieves $R^2 > 0.98$.

7 Experiments

7.1 Main Results

We compare DLCM against a parameter-matched baseline that follows the LLaMA [18] architecture. Both models are trained from scratch on our proprietary dataset, using the same global batch size, learning rate, and sequence length as reported in the LLaMA paper. Each model is trained on 1T tokens. Results on 12 standard zero-shot benchmarks are summarized in Table 2.

Our model follows an encoder–compressor–decoder architecture with learned concept circulation, explicitly redistributing computation from uniform token-level processing to adaptive concept-level reasoning. As a consequence, we do not expect uniform gains across all benchmarks. Instead, performance differences directly reflect the architectural bias induced by semantic compression and boundary-aware compute allocation.

Overall, DLCM achieves an average accuracy of **43.92%**, surpassing the baseline score of 41.23% by **+2.69%**. However, these gains are highly non-uniform across tasks, revealing a clear separation between reasoning-dominant benchmarks and those that rely on fine-grained token-level alignment.

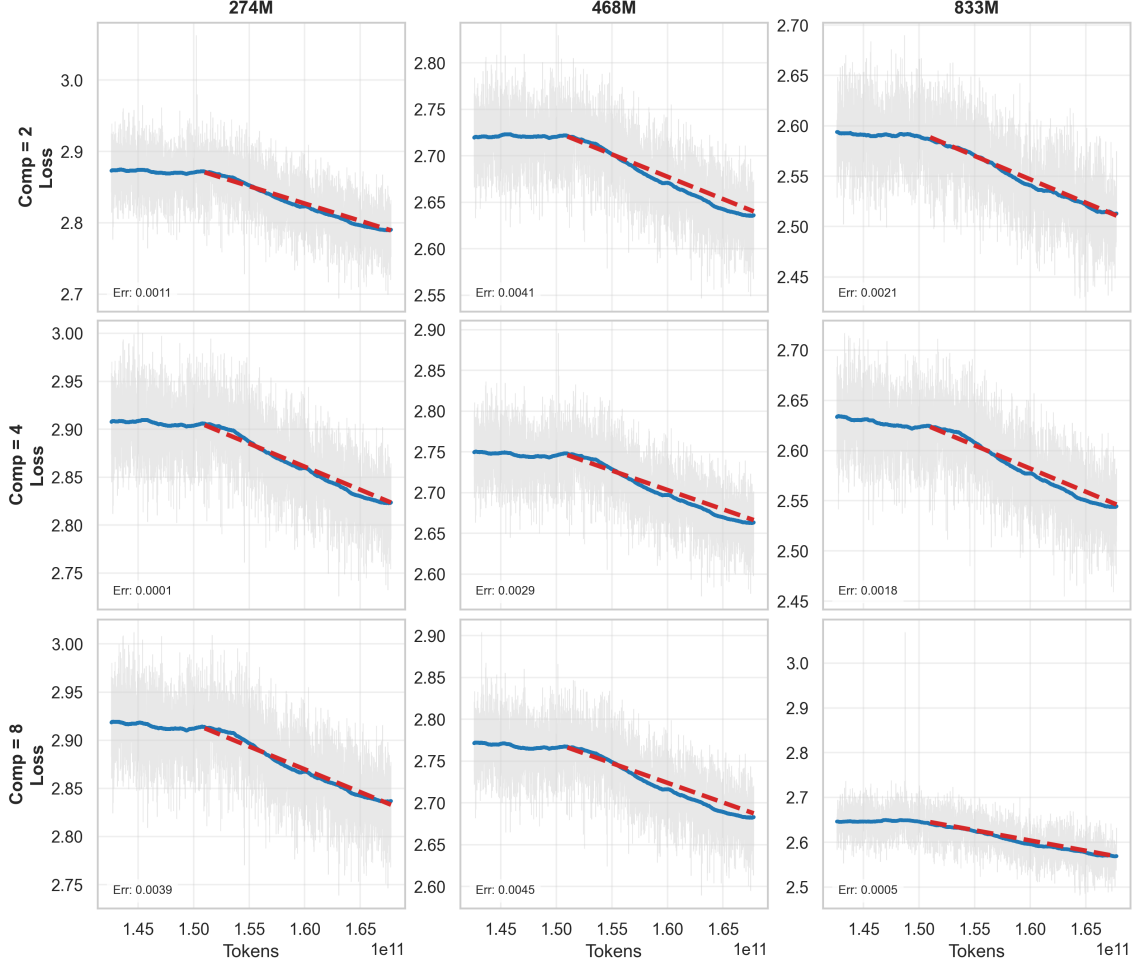


Figure 5 Decay-phase fit. Simplified fit on the final portion of training tokens, validating that our WSD scaling law accurately captures late-stage behaviour with $R^2 = 0.93$.

Reasoning-Dominant Tasks. We observe consistent and often substantial improvements on benchmarks that emphasize multi-step reasoning, hypothesis selection, and implicit commonsense inference. Notable gains are achieved on **CommonsenseQA (+1.64%)**, **HellaSwag (+0.67%)**, **OpenBookQA (+3.00%)**, **PIQA (+2.42%)**, and both **ARC Easy (+2.61%)** and **ARC Challenge (+1.77%)**. These tasks are characterized by non-uniform information density, where prediction difficulty concentrates around semantic transitions rather than being evenly distributed across tokens. By compressing locally predictable spans and allocating the majority of model capacity to a high-dimensional concept backbone, DLCM focuses computation on structurally salient regions. This behavior is consistent with our loss distribution analysis in Section 7.2, which shows systematic loss reduction near concept boundaries.

Granularity-Sensitive Text Understanding. In contrast, we observe mild regressions on **BoolQ (-1.47%)** and **RACE (-0.72%)**. These benchmarks depend heavily on fine-grained sentence-level entailment, polarity resolution, and subtle lexical cues. The encoder-compress-decode paradigm inevitably reduces token-level granularity within concept interiors, which can obscure micro-level distinctions required for such tasks. Importantly, this degradation is localized rather than uniform: while boundary tokens are modeled more accurately, mid-concept positions may trade off fine-grained precision for improved global coherence. This trade-off manifests as the U-shaped loss profile observed in our mechanistic analysis.

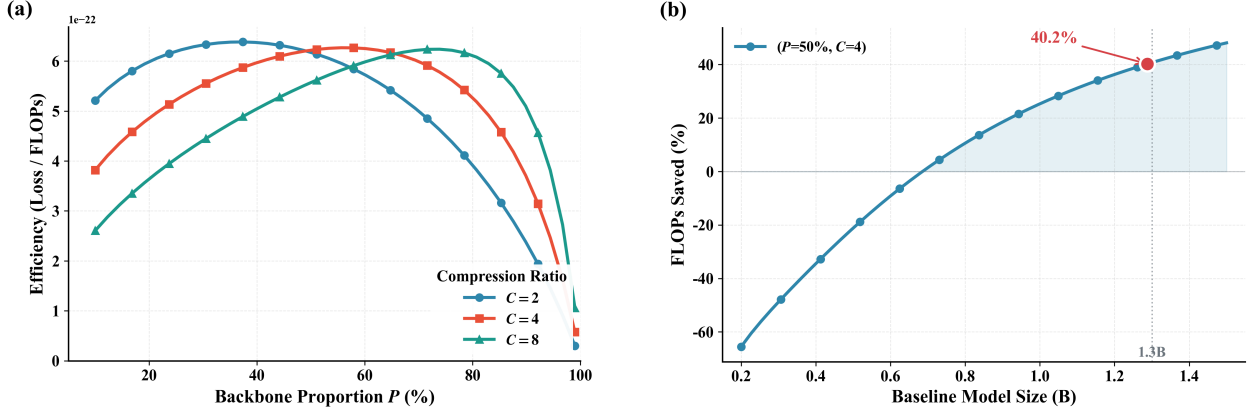


Figure 6 Efficiency analysis of the DLCM architecture. (a) Architectural efficiency (Loss/FLOPs) across backbone proportions P for different compression ratios R . (b) FLOPs savings compared to baseline models of varying sizes, with DLCM configured at $P = 60\%$, $R = 4$.

Knowledge and Multilingual Benchmarks. For encyclopedic knowledge evaluation, we observe mixed behavior. While **C-Eval (+1.71%)** benefits from adaptive segmentation enabled by the Global Parser, slight regressions appear on **MMLU (-0.30%)** and **CMMLU (-0.24%)**. These datasets reward relatively uniform factual recall across tokens, leaving less opportunity for boundary-aware compute reallocation. This result further supports our central claim: DLCM is structurally optimized for reasoning under non-uniform information density, rather than uniform memorization-heavy retrieval.

Architecture and Parameter Efficiency. Although DLCM contains nearly $2\times$ the total parameters of the baseline model (2.3B vs. 1.3B), this increase is deliberately concentrated in the concept-level backbone ($d_{\text{concept}} = 3072$). Because this backbone operates on a sequence compressed by $4\times$, the effective FLOPs per inference step remain comparable to the smaller baseline. This validates our core design principle: shifting computation from redundant token-level processing to dense concept-level reasoning enables substantially larger effective capacity without incurring proportional inference cost.

7.2 Analysis: Compute Allocation in Concept-Based Models

7.2.1 Experimental Setup

To isolate the impact of concept-based compression on model behavior, we conduct a controlled comparison between our proposed concept model and a standard Transformer baseline. Both models utilize the same backbone architecture (1.3B parameters) and were trained on an identical subset of 100B tokens from our pretraining corpus. This ensures that any observed differences in loss distribution are attributable solely to the compression mechanism and architectural changes, rather than discrepancies in training data or compute budget.

7.2.2 Loss Distribution Analysis

To understand how the model allocates computational resources, we evaluate the loss distribution across relative positions within concepts. We randomly selected 600 samples from the validation set and aligned the token-level losses based on their position within a segmented concept (e.g., the i -th token of a concept).

Figure 7 illustrates the average loss at the first 20 positions within each concept. The top panel compares raw loss values, while the bottom panel visualizes the differential: $\Delta L = L_{\text{concept}} - L_{\text{baseline}}$. Here, **green bars** ($\Delta L < 0$) indicate the concept model outperforms the baseline, while **red bars** ($\Delta L > 0$) indicate degradation.

The results reveal a distinct "U-shaped" improvement pattern that reflects the model's resource reallocation strategy:

Table 2 Performance Comparison: DLCM vs. Baseline. Zero-shot accuracy (%) categorized by task type. Improvements are shown in green and regressions in red.

| Task / Category | DLCM (Ours) | Baseline | Diff. |
|--|--------------|--------------|-------|
| <i>Multi-choice General Knowledge / Common Sense</i> | | | |
| Commonsense QA | 21.38 | 19.74 | +1.64 |
| HellaSwag | 46.66 | 45.99 | +0.67 |
| Winogrande | 57.22 | 56.20 | +1.02 |
| OpenBookQA | 26.80 | 23.80 | +3.00 |
| PIQA | 75.52 | 73.10 | +2.42 |
| ARC Challenge | 34.81 | 33.04 | +1.77 |
| ARC Easy | 69.91 | 67.30 | +2.61 |
| MMLU | 25.40 | 25.70 | -0.30 |
| <i>Multi-choice Text Understanding</i> | | | |
| BoolQ | 62.54 | 64.01 | -1.47 |
| RACE | 35.31 | 36.03 | -0.72 |
| <i>Culture / Multilingual Knowledge</i> | | | |
| C-Eval | 26.08 | 24.37 | +1.71 |
| CMMLU | 25.23 | 25.47 | -0.24 |
| Average | 43.92 | 41.23 | +2.69 |

1. **Boundary Proficiency (Positions 0--2 & 16+):** Consistent with the observation that concept models excel at initial and late positions (indicated by green bars), the architecture effectively captures the transition semantics. By explicitly modeling concept boundaries, the model reduces ambiguity at the start and end of semantic units, outperforming the baseline which treats these tokens uniformly.
2. **Internal Complexity (Mid-positions):** In the middle of a concept (approx. positions 4–15), we observe a shift. While the baseline model often struggles here (higher absolute loss), the concept model’s performance is mixed. The presence of red bars in certain mid-concept regions suggests that the compression mechanism forces the model to trade off some fine-grained token-level precision to maintain higher-level semantic coherence.

This reallocation aligns with our hypothesis: the concept model sacrifices uniform token-level predictability (resulting in minor degradation at specific internal positions) to gain superior performance at semantic boundaries and structurally critical tokens. This strategic trade-off allows the model to "spend" its capacity on maintaining global coherence, explaining the downstream improvements despite non-uniform loss reduction.

8 Ablation Studies

8.1 Analysis: End-to-End Discrete Boundary Learning vs. Decoupled Segmentation

This experiment compares two boundary prediction mechanisms for sequence compression: a learned neural predictor with compression rate regularization (Section 3.3.3) and a rule-based predictor using cosine similarity. Starting with sequences of length $L = 8192$, we track the average compressed length during training.

Figure 8 reveals starkly different behaviors. The **learned predictor (red)** exhibits severe instability: after initial compression to ~ 2000 tokens, the compressed length steadily increases, eventually stabilizing at ~ 4300 tokens ($1.9\times$ compression). This "creep-up" indicates the model progressively learns to compress less over time. In contrast, the **rule-based predictor (purple)** demonstrates exceptional stability, rapidly converging to ~ 2000 tokens ($4\times$ compression) and maintaining this level consistently throughout training.

Table 3 Architecture Configuration Details. A unified view of the parameter settings for Baseline (LLaMA-1.3B) and DLCM (2.3B). Values are presented as *Baseline* / *Ours*.

| General Settings | | Dimension Settings | |
|---------------------|----------------------|--------------------------------------|------------------|
| <i>Metric</i> | <i>Value</i> | <i>Metric</i> | <i>Value</i> |
| Model Type | Trans. / DLCM | Hidden Size (d_{token}) | 1,536 |
| Total Params | 1.3B / 2.3B | Main Hidden (d_{concept}) | – / 3,072 |
| Vocab Size | 128,815 | Interm. Size (Self) | 4,096 / 6,144 |
| Max Pos Emb | 8k / 8k | Interm. Size (Cross) | – / 6,144 |
| Activation | Swish | | |
| Layer Configuration | | Attention Configuration | |
| Total Layers | 32 | Attn Heads | 24 / 24 |
| Encoder Layers | – / 10 | Backbone Heads | – / 48 |
| Backbone Layers | – / 16 | KV Heads | 24 / 12 |
| Decoder Layers | – / 6 | Backbone KV | – / 24 |

The learned predictor’s instability stems from conflicting optimization objectives. Despite the compression rate regularization term \mathcal{L}_{aux} designed to maintain the target compression ratio R , the primary cross-entropy (CE) loss creates much stronger gradients that penalize information loss and discourage compression:

$$\nabla_{\theta} \mathcal{L}_{\text{total}} = \underbrace{\nabla_{\theta} \mathcal{L}_{\text{CE}}}_{\text{anti-compression}} + \lambda \underbrace{\nabla_{\theta} \mathcal{L}_{\text{aux}}}_{\text{pro-compression}} \quad (24)$$

Since $\|\nabla_{\theta} \mathcal{L}_{\text{CE}}\| \gg \lambda \|\nabla_{\theta} \mathcal{L}_{\text{aux}}\|$, the CE loss eventually dominates, forcing the predictor to reduce segmentation despite the regularization term.

The rule-based predictor avoids this conflict through a fixed decision rule: $p_t = \frac{1 - \cos(h_t, h_{t+1})}{2}$, with boundaries inserted when $p_t > \tau$. While the representations h_t are learned, the segmentation rule itself is not optimized by the CE loss. This decoupling prevents the task loss from undermining the compression mechanism, ensuring stable and controllable compression ratios through the threshold parameter τ .

8.2 Global Regularization via Gradient Accumulation

To further stabilize the learned boundary predictor, we investigate an alternative regularization strategy: computing the compression ratio loss over accumulated training examples rather than individual sequences (Section 3.3.3). This **global regularization** approach computes boundary statistics F_{global} and G_{global} across all tokens in K micro-batches.

We train two 2.3B parameter models for 1T tokens with a target compression ratio of $R = 2$: one with per-sequence regularization ("Normal") and one with global regularization ("Global Parser"). Table 4 presents the downstream performance and the actual realized compression ratios.

The global regularization approach achieves consistently better performance across most tasks (5 out of 6). Crucially, as shown in the bottom row of Table 4, the Global Parser maintains a realized compression ratio (~ 3.9) much closer to the target (4.0) compared to the Normal formulation, which tends to degrade towards lower compression.

The key insight is that enforcing a fixed compression ratio per sequence is overly restrictive. Real-world data exhibits varying information density. By relaxing the constraint to operate at the batch level, the global regularization allows the model to learn adaptive behavior—compressing repetitive code more aggressively while preserving dense technical text—effectively allocating the compression "budget" where it matters most.

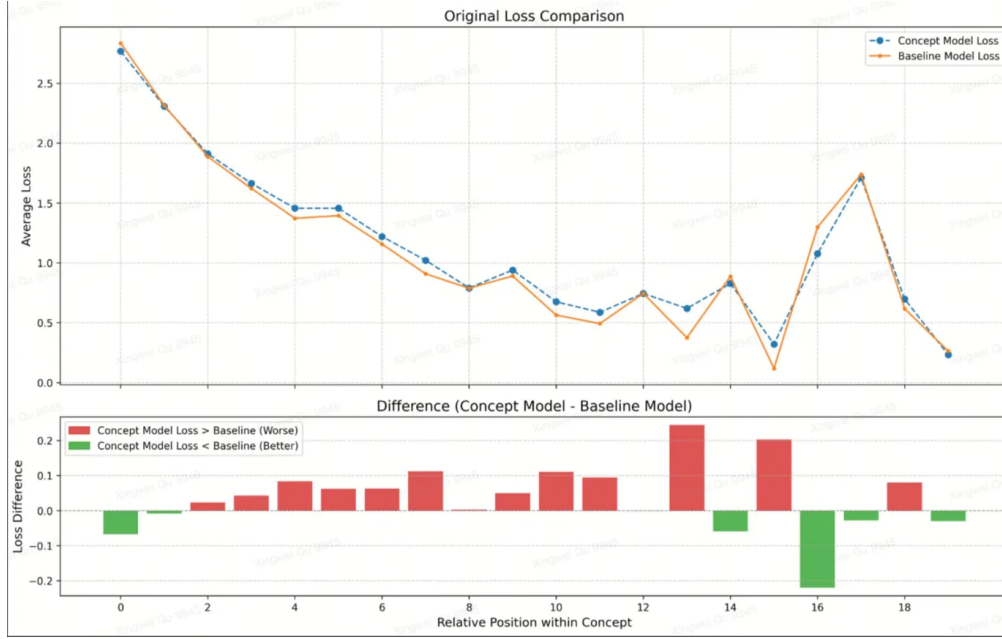


Figure 7 Top: Average loss comparison between concept model (blue) and baseline model (orange) across relative positions within concepts. Bottom: Loss difference (Concept - Baseline), where **green indicates improvement** (lower loss) and red indicates degradation.

8.3 Content-Adaptive Compression Benefits

To verify the adaptive behavior enabled by global regularization, we analyzed the segmentation granularity across different domains. Table 5 presents empirical measurements of the average tokens per concept.

The data reveals significant variation in compression density across content types. For instance, at the $8\times$ target, Technical English retains significantly more tokens per concept (10.58) compared to Technical Chinese (6.09) or Code (6.14).

While the precise ranking of "optimal" length varies across compression targets (as noted in the fluctuation between content types), the **existence of this variation** is the critical finding. It confirms that the global regularization mechanism successfully decouples the compression objective from rigid per-sequence constraints. The model is not forcing a uniform segment length; instead, it adapts the granularity based on the inherent

Table 4 Ablation Study: Global Parser vs. Normal. Performance comparison on downstream tasks. Both models aim for a target compression ratio of $R = 4$. The **Global Parser** achieves a realized ratio much closer to the target while consistently improving accuracy on most tasks.

| Task | Global Parser | Normal | Metric |
|-------------------------|---------------|---------------|-----------------|
| ARC Challenge | 0.3038 | 0.2858 | Acc |
| ARC Easy | 0.6296 | 0.6242 | Acc |
| Commonsense QA | 0.2457 | 0.2228 | Acc |
| HellaSwag | 0.3507 | 0.3499 | Acc |
| OpenBookQA | 0.3220 | 0.3280 | Acc |
| PIQA | 0.6806 | 0.6785 | Acc |
| Avg. Improvement | +2.1% | — | — |
| Realized Ratio | 3.92 | 3.15 | (Target $R=4$) |

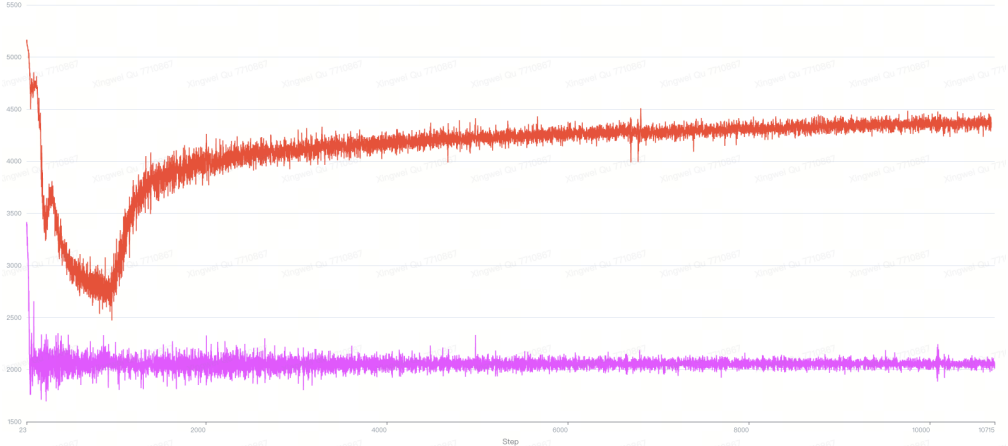


Figure 8 Average compressed sequence length over training steps. **Red:** Learned Boundary Predictor. **Purple:** Rule-Based Predictor. The x-axis represents training steps, and the y-axis represents the average number of tokens post-compression.

| Content Type | Target 8× | Target 4× | Target 2× |
|-------------------|-----------|-----------|-----------|
| Casual English | 7.47 | 3.53 | 1.76 |
| Casual Chinese | 8.38 | 4.36 | 1.76 |
| Technical English | 10.58 | 3.85 | 1.92 |
| Technical Chinese | 6.09 | 3.27 | 1.76 |
| Code | 6.14 | 3.66 | 1.98 |
| Math/Science | 7.42 | 4.41 | 1.91 |

Table 5 Average tokens per concept across content types and compression ratios. Values represent the actual granularity achieved for each target compression setting.

semantic density of the content. Code and structured text tend to be compressed into shorter, syntactic units, whereas dense prose is preserved in longer semantic chunks. This adaptivity—regardless of the specific order—allows the model to maximize information retention within the global compression budget.

9 Conclusion

We presented **Dynamic Large Concept Models (DLCM)**, a hierarchical language modeling framework that challenges the token-uniform computation paradigm underlying modern LLMs. By learning semantic boundaries from latent representations and shifting computation from tokens to variable-length concepts, DLCM enables reasoning to occur in a compact, semantically aligned space rather than repeatedly at the token level.

Beyond the architectural design, we showed that hierarchical compression necessitates new theoretical and optimization tools. We introduced a compression-aware scaling law that clarifies how compute should be allocated between token processing and concept-level reasoning under fixed FLOPs, and developed a decoupled μP parametrization that enables stable training and zero-shot hyperparameter transfer in heterogeneous architectures. Empirically, DLCM achieves consistent gains on reasoning-intensive benchmarks while reducing redundant computation, demonstrating a favorable accuracy–efficiency trade-off.

More broadly, our results suggest that scaling language models is not solely a matter of increasing parameters or data, but also of reconsidering where computation is performed. We believe concept-level latent reasoning offers a promising direction for building more efficient and more reasoning-capable language models, and opens avenues for future work on adaptive abstraction, planning, and multi-level reasoning in large-scale neural systems.

Contributions

Leading Authors

Xingwei Qu, Shaowen Wang, Zihao Huang, Ge Zhang

Leading Author Contributions

Xingwei Qu: Conducts fundamental ablation studies and is responsible for the majority of the engineering implementation.

Shaowen Wang: Implements and advances the MuP (Maximal Update Parametrization) for hyperparameter tuning.

Zihao Huang: Implements noise based boundary prediction tricks.

Ge Zhang: Proposes the original idea and develops the demo prototype. Identifies and provides the key technique for the Global Parser.

Core Contributors

Kai Hua: Designs and constructs the training data entirely from open-source data.

Fan Yin: Contributes to the SGLang implementation and optimization.

Rui-Jie Zhu: Resolves bugs related to DLCM and proposes the strategy of replacing the encoder with DLCM.

Jundong Zhou & Qiyang Min: Contributes to the project development and implementation.

Other Contributors

Zihao Wang, Yizhi Li, Tianyu Zhang, He Xing, Zheng Zhang, Yuxuan Song, Tianyu Zheng, Zhiyuan Zeng

Corresponding Authors

Chenghua Lin, Ge Zhang, Wenhao Huang

References

- [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [2] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359, 2022.
- [3] Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Łukasz Kaiser. Universal transformers. *arXiv preprint arXiv:1807.03819*, 2018.
- [4] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, Rodolphe Jenatton, Lucas Beyer, Michael Tschannen, Anurag Arnab, Xiao Wang, Carlos Riquelme, Matthias Minderer, Joan Puigcerver, Utku Evci, Manoj Kumar, Sjoerd van Steenkiste, Gamaleldin F. Elsayed, Aravindh Mahendran, Fisher Yu, and Neil Abnar, Samira andXHoulsby. Scaling vision transformers to 22 billion parameters. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 7480–7512. PMLR, 2023. URL <https://proceedings.mlr.press/v202/dehghani23a.html>.
- [5] Xinrun Du, Zhouliang Yu, Songyang Gao, Ding Pan, Yuyang Cheng, Ziyang Ma, Ruibin Yuan, Xingwei Qu, Jiaheng Liu, Tianyu Zheng, et al. Chinese tiny llm: Pretraining a chinese-centric large language model. *arXiv preprint arXiv:2404.04167*, 2024.
- [6] Sirui Hao, Yeyun Shen, Weizhu Chen, Dongyan Zhang, and Ming Zhou. Training large language models to reason in a continuous latent space. *arXiv preprint arXiv:2412.06769*, 2024.
- [7] Alex Henry, Prudhvi Raj Dachapally, Shubham Shantaram Pawar, and Yuxuan Chen. Query-Key Normalization for Transformers. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4246–4253. Association for Computational Linguistics, 2020. doi: 10.18653/v1/2020.findings-emnlp.379. URL <https://aclanthology.org/2020.findings-emnlp.379>.
- [8] 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. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- [9] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, and et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022. URL <https://arxiv.org/abs/2203.15556>.
- [10] Siming Huang, Tianhao Cheng, Jason Klein Liu, Jiaran Hao, Liuyihan Song, Yang Xu, J. Yang, J. H. Liu, Chenchen Zhang, Linzheng Chai, Ruifeng Yuan, Zhaoxiang Zhang, Jie Fu, Qian Liu, Ge Zhang, Zili Wang, Yuan Qi, Yinghui Xu, and Wei Chu. Opencoder: The open cookbook for top-tier code large language models. 2024. URL <https://arxiv.org/pdf/2411.04905>.
- [11] Sukjun Hwang, Brandon Wang, and Albert Gu. Dynamic chunking for end-to-end hierarchical sequence modeling, 2025. URL <https://arxiv.org/abs/2507.07955>.
- [12] Albert Q Jiang, Alexandre Sablayrolles, Arthur Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra S Chaplot, Diego de las Casas, Emma B Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- [13] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [14] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- [15] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarczyk, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. In *International Conference on Learning Representations*, 2017.
- [16] Dan Su, Kezhi Kong, Ying Lin, Joseph Jennings, Brandon Norick, Markus Kliegl, Mostofa Patwary, Mohammad Shoyebi, and Bryan Catanzaro. Nemotron-cc: Transforming common crawl into a refined long-horizon pretraining dataset. *arXiv preprint arXiv:2412.02595*, 2024.

- [17] LCM team, Loïc Barrault, Paul-Ambroise Duquenne, Maha Elbayad, Artyom Kozhevnikov, et al. Large concept models: Language modeling in a sentence representation space. [arXiv preprint arXiv:2412.08821](#), 2024.
- [18] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. [arXiv preprint arXiv:2302.13971](#), 2023.
- [19] Greg Yang, Edward J Hu, Igor Babuschkin, Szymon Sidor, Xiaodong Liu, David Farhi, Nick Ryder, Jakub Pachocki, Weizhu Chen, and Jianfeng Gao. Tensor programs v: Tuning large neural networks via zero-shot hyperparameter transfer. [arXiv preprint arXiv:2203.03466](#), 2022.
- [20] Fan Zhou, Zengzhi Wang, Nikhil Ranjan, Zhoujun Cheng, Liping Tang, Guowei He, Zhengzhong Liu, and Eric P. Xing. Megamath: Pushing the limits of open math corpora. [arXiv preprint arXiv:2504.02807](#), 2025. Preprint.

Appendix

A Appendix: Segmentation Examples at Different Compression Ratios

We provide representative examples of how boundary prediction behaves under different compression ratios across three content types: casual English text, Python code, and mathematical exposition.

A.1 Casual English Text

Original (90 tokens):

So I've been trying to perfect my morning coffee routine lately. It's funny how something so simple can have so many variables. I started with a basic drip machine, which was fine, but a bit boring. Then I went down the French press rabbit hole – way more flavor, but the cleanup is a real hassle, you know?

Compression 8× (11 segments):

So I | 've been trying to perfect my morning coffee routine lately. | It's funny how something | so simple can have so many variables. | I started with | a basic drip machine, which | was fine, but a | bit boring. | Then I went down the French press rabbit hole – | way more flavor, but the cleanup is a | real hassle, you know?

Compression 4× (35 segments):

So I | 've been trying | to perfect | my morning coffee routine lately. | It's | funny how something | so simple can | have so many variables. | I started | with | a basic drip machine, which was fine, but | a bit boring.

Compression 2× (56 segments):

So I | 've | been trying to | perfect | my | morning coffee routine | lately. | It | 's funny how | something so | simple | can | have so many | variables.

A.2 Python Code

Original (87 tokens):

```
import torch
from torch.utils.data import Dataset, DataLoader
class SimpleTextDataset(Dataset):
    """A simple dataset for loading text data."""
    def __init__(self, texts, tokenizer, max_length=128):
        self.texts = texts
        self.tokenizer = tokenizer
```

Compression 8× (12 segments):

```
import torch | from torch.utils.data import Dataset, DataLoader |
class SimpleTextDataset(Dataset): |
    """A simple dataset for loading text data.""" |
    def __init__(self, texts, tokenizer, max_length=128): |
self.texts = | texts | self.tokenizer = | tokenizer
```

Compression 4× (29 segments):

```
import torch | from torch.utils.data | import Dataset | , DataLoader |
class Simple | TextDataset(Dataset): | """ | A | simple dataset for |
def | __init__(self, texts, tokenizer, max_length= | 128): |
self.texts = | texts
```

Compression 2× (58 segments):


```
import | torch | from torch.utils.data import | Dataset, | DataLoader |
class | Simple | TextDataset(Dataset): | "" | A simple | dataset for |
def | __init__(self, | texts, | tokenizer | , | max_length= | 128): |
self | .texts | = | texts
```

A.3 Mathematical Text

Original (95 tokens):

Euler’s formula is a mathematical formula in complex analysis that establishes the fundamental relationship between the trigonometric functions and the complex exponential function. The formula states that for any real number x , $e^{ix} = \cos(x) + i \sin(x)$, where e is the base of the natural logarithm, i is the imaginary unit.

Compression $8\times$ (13 segments):

Euler’s formula is a mathematical formula in complex analysis that establishes the fundamental relationship between the trigonometric functions and the complex exponential function. The formula states that for any real number x , $e^{ix} = \cos(x) + i \sin(x)$, where e is the base of the natural logarithm, i is the imaginary unit

Compression $4\times$ (32 segments):

Euler’s formula is a mathematical formula in complex analysis that establishes the fundamental relationship between the trigonometric functions and the complex exponential function. The formula states that for any real number x , $e^{ix} = \cos(x) + i \sin(x)$, where e is the base of the natural logarithm

Compression $2\times$ (61 segments):

Euler’s formula is a mathematical formula in complex analysis that establishes the fundamental relationship between the trigonometric functions and the complex exponential function. The formula states that for any real number x , $e^{ix} = \cos(x) + i \sin(x)$,

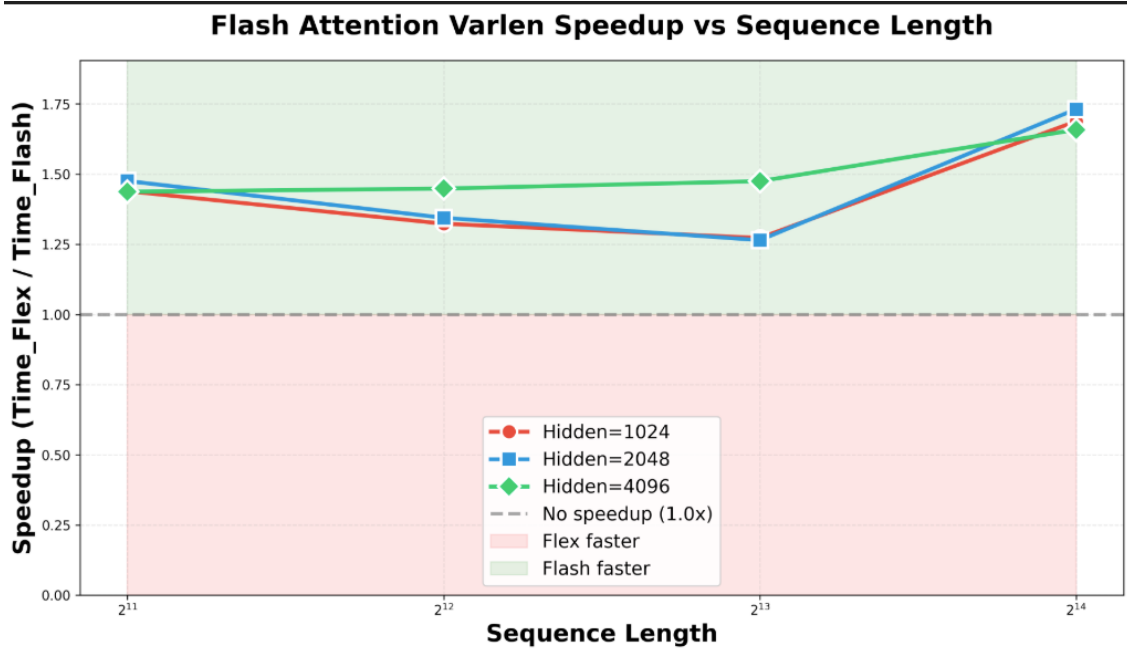


Figure 9 Flash Attention Varlen speedup ($T_{\text{flex}}/T_{\text{flash}}$) vs. Sequence Length. The plot visualizes the data from Table 6, highlighting the performance trend across different scales and hidden sizes.

Table 6 Performance comparison (Batch=1, Heads=32, Interval=6)

| Seq Length | Hidden Size | Flex (ms) | Flash Varlen (ms) | Speedup |
|------------|-------------|-----------|-------------------|---------------|
| 2048 | 1024 | 32.35 | 22.48 | 1.44 × |
| 2048 | 2048 | 33.31 | 22.58 | 1.48 × |
| 2048 | 4096 | 32.42 | 22.56 | 1.44 × |
| 4096 | 1024 | 59.75 | 45.15 | 1.32 × |
| 4096 | 2048 | 60.72 | 45.17 | 1.34 × |
| 4096 | 4096 | 65.88 | 45.48 | 1.45 × |
| 8192 | 1024 | 116.35 | 91.42 | 1.27 × |
| 8192 | 2048 | 114.65 | 90.66 | 1.26 × |
| 8192 | 4096 | 142.75 | 96.79 | 1.47 × |
| 16384 | 1024 | 314.35 | 186.21 | 1.69 × |
| 16384 | 2048 | 323.53 | 186.83 | 1.73 × |
| 16384 | 4096 | 315.69 | 190.38 | 1.66 × |