
Chain-of-Model Learning for Language Model

Xiaohua Wang^{1,2,*} Kaitao Song^{2,*,†} Xu Tan^{*} Huiqiang Jiang² Chengruidong Zhang²
Yongliang Shen³ Cen Lu^{2,5,6} Zihao Li^{2,7} Zifan Song^{2,8}
Caihua Shan² Yansen Wang² Kan Ren^{2,4} Xiaoqing Zheng¹
Tao Qin² Yuqing Yang² Dongsheng Li² Lili Qiu²

Fudan University¹ Microsoft Research² Zhejiang University³ ShanghaiTech University⁴
Idiap Research Institute⁵ EPFL⁶ UIUC⁷ Tongji University⁸

Abstract

In this paper, we propose a novel learning paradigm, termed “*Chain-of-Model*” (CoM), which incorporates the causal relationship into the hidden states of each layer as a chain style, thereby introducing great scaling efficiency in model training and inference flexibility in deployment. We introduce the concept of “*Chain-of-Representation*” (CoR), which formulates the hidden states at each layer as a combination of multiple sub-representations (i.e., chains) at the hidden dimension level. In each layer, each chain from the output representations can only view all of its preceding chains in the input representations. Consequently, the model built upon CoM framework can progressively scale up the model size by increasing the chains based on the previous models (i.e., chains), and offer multiple sub-models at varying sizes for elastic inference by using different chain numbers. Based on this principle, we devise *Chain-of-Language-Model* (CoLM), which incorporates the idea of CoM into each layer of Transformer architecture. Based on CoLM, we further introduce CoLM-Air by introducing a *KV sharing* mechanism, that computes all keys and values within the first chain and then shares across all chains. This design demonstrates additional extensibility, such as enabling seamless LM switching, prefilling acceleration and so on. Experimental results demonstrate our CoLM family can achieve comparable performance to the standard Transformer, while simultaneously enabling greater flexibility, such as progressive scaling to improve training efficiency and offer multiple varying model sizes for elastic inference, paving a new way toward building language models. Our code will be released in the future at: <https://github.com/microsoft/CoLM>.

1 Introduction

With the advent of large language models (LLMs) [1–6], scaling up Transformer architecture [7–9] has been considered as the de facto approach to revolutionize existing AI landscape and achieve state-of-the-art performance across a wide range of different tasks. Therefore, there has a growing trend in both industry and academia to explore strategies for scaling up Transformer models. Under these circumstances, the parameter size of LLMs has been an exponential growth, scaling from the billion to the trillion level. As a result, its explosive parameters also impose extremely expensive burden on training and cannot offer varying inference usage for different deployment environments. In light of this increasing trend of scaling laws, how to develop and harness LLMs effectively to address user instructions in various scenarios has been an open and critical challenge among the whole community.

* The first three authors have equal contributions.

† Corresponding author: Kaitao Song (stillkeeptry@outlook.com)

Motivated by elastic inference [10, 11] and continual training [12, 13], we claim that existing scaling strategies for LLM architectures still remain some inherent issues:

- Unlike human intelligence, which acquires new knowledge incrementally, existing scaling policies cannot retain existing scales and always requires training from scratch, leading to inefficiencies.
- Existing LLM architectures (e.g., Dense or MoE [14]) always activate a fixed scale of parameters, lacking mechanisms to dynamically adapt problem-solving capabilities.

In this paper, we propose a new concept, termed **“Chain-of-Representation”** (CoR), that generalizes the scope of representation paradigm to a broader range. Specifically, we observe that any representation can always be viewed as a combination of multiple sub-representations at the hidden dimension level. Therefore, we define this composition as Chain-of-Representation, that each sub-representation corresponds to one chain. Based on this definition, by using different number of preceding chains, its corresponding features can be used to encode different knowledge (we call it “scale”), just as shown in Figure 1. Therefore, how to build the connections among CoR features to ensure feature transformation across each scale has been a critical step.

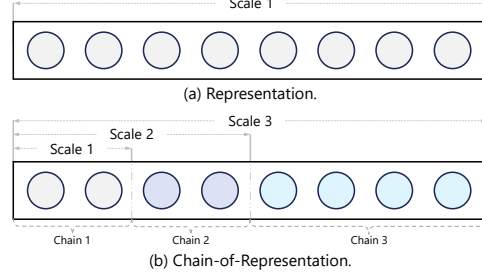


Figure 1: An example of CoR with 3 chains.

To this end, we introduce **“Chain-of-Model”** (CoM), a novel learning paradigm for modeling CoR features. Its core idea is to incorporate causal dependencies across different scales, ensuring that each scale can only use information from its preceding scales at the representation level. To fulfill this, we propose **“Chain-of-Layer”** (CoL), to reformulate current network layers based on CoR features. Specifically, we require each output chain i is only conditioned on the input chains from 1 to i . Figure 2 illustrates an example of applying CoL into Linear layer. CoL is capable of three important characteristics: *Generality*, *Causality* and *Compositionality*. That means, 1) Any layer can be viewed as CoL (chain is 1); 2) To obtain the features at scale i , CoL allows us to only activate the parameters of its correlated chains; 3) If two layers both conform to the CoL settings, their composition also adheres to the CoL property. Following this principle, we can stack multiple CoL layers to generalize it from the layer level to the model level, i.e., Chain-of-Model (CoM). The scope of CoM can be generalized to any ML models as all models can be considered as a special case of CoM that chain is set as 1. Besides, the output features of CoM also conform to the characteristic of CoR, CoM can offer multiple sub-models at different scales by using different number of chains. Based on these characteristics, CoM can exhibit greater flexibility to build foundation models.

On the basis of CoM framework, we re-formulate the architecture of language models by applying the idea of CoL into each layer within Transformer, and term it as **“Chain-of-Language-Model”** (CoLM). CoLM can integrate multi-scale training within single forward propagation, to offer multiple sub-models for elastic inference. It also enables us to scale up language models via chain expansion (i.e., increase the number of chains), while simultaneously preserving the capability to use previous scales. Moreover, based on the CoL criteria, we further introduce a *KV sharing* mechanism in attention module that requires all keys and values to be computed within the first chain, and term it as CoLM-Air. Based on this mechanism, CoLM-Air offers greater extensibility and flexibility, including allowing us to seamless switching any different scales of LLMs without any additional re-computations (e.g., keys and values), accelerating prefilling within the first chain and so on. Experimental results on multiple benchmarks also demonstrate our CoLM family can achieve comparable performance while both exhibit better extensibility and flexibility. Overall, the contributions of our paper can be summarized as below:

- We point out the inherent issues of existing LLM scaling approaches, and formulate the concept of chain-of-representation (CoR), that integrates multi-scale capabilities within single representation.
- To model the relationships among CoR features, we introduce chain-of-model (CoM) learning and further extend it into language model, termed CoLM. CoLM is able to provide multiple sub-models at varying scales within a unified foundation model, enabling more flexible and scalable extensions.
- Experimental results on multiple benchmarks demonstrate that our CoLM family can achieve comparable performance to existing LLM frameworks, while simultaneously exhibit better extension extensibility and flexibility in different scenarios (e.g., prefilling, tuning and so on).

2 Chain-of-Model Learning

In this section, we first introduce the concept of “*Chain-of-Representation*” (CoR), which generalizes the paradigms of existing representation to a broader scope. After that, we introduce “*Chain-of-Layer*” (CoL) and “*Chain-of-Model*” (CoM), to model CoR from the layer level to the model level.

2.1 Representation

Definition 1 (*Chain-of-Representation*) For representation $x \in \mathbb{R}^D$, it can always be equivalently formulated as a concatenation of n sub-representations, denoted as $\xi(x, n) = \{x_1, \dots, x_n\}$, where $x_i \in \mathbb{R}^{d_i}$ and $\sum_{i=1}^n d_i = D$. We term this as chain-of-representation (CoR) of x .

According to Definition 1, each chain corresponds to each sub-representation (i.e., x_i) in CoR. By activating the first few chains (i.e., $x_{\leq i}$), it can be used to encode information of scale _{i} , with a dimension of $\text{sum}(d_{\leq i})$. Hence, CoR allows us to encode n different scales within one representation. If $n = 1$, CoR is identical to the original representation. Figure 1 illustrates the CoR.

2.2 Layer

However, a challenge arises is how to design layers to build connections between CoR inputs and CoR outputs, enabling multi-scale feature transformation while both preserving the output features to follow the criteria of CoR in Definition 1. Therefore, we need to maintain each scale can only utilize information from all its preceding scales, and introduce *Chain-of-Layer* to incorporate causal relationships into hidden states of CoR, as below:

Definition 2 (*Chain-of-Layer*) For a layer $y = f_\theta(x)$, where its input x and output y can both be expressed as the chain-of-representation $\xi(x, n)$ and $\xi(y, n)$. We term it as a Chain-of-Layer (CoL) where $f_\theta(\cdot)$ satisfies that each y_i is only conditioned on $x_{\leq i}$.

CoL is capable of three fundamental properties – *generality*, *causality* and *compositionality*:

Corollary 2.1 (*Generality*) Any layer can be viewed as a case of chain-of-layer when $n = 1$.

If $n = 1$, CoL can be considered as the standard network layer to process the conventional representation. In other words, any network layer can be generalized to the form of CoL. Therefore, we can gradually increase layer dimension by introducing additional chains upon the existing chains.

Corollary 2.2 (*Causality*) If a layer $y = f(x)$ satisfies the property of chain-of-layer, that means its weights θ can be partitioned into n independent parts $\{\theta_1, \dots, \theta_n\}$, where each θ_i is used to compute y_i based on $x_{\leq i}$. Therefore, y allows the causal computation, i.e., to obtain feature y_i , we only need to compute $\theta_{\leq i}$ based on $x_{\leq i}$.

In the CoL design, each scale can aggregate information of all its previous scales to increase scalability and effectively avoid catastrophic forgetting [15]. Therefore, to obtain information of Scale _{i} , this design allows us to only compute its preceding chains (i.e., $\theta_{\leq i}$), without calculating the whole weights. As a result, CoL can offer multi-scale information within one representation and allow us to dynamically calculate features at varying scales.

Corollary 2.3 (*Compositionality*) Assume we have two layers $y = f_1(x)$ and $z = f_2(y)$, where $x = \xi(x, n)$, $y = \xi(y, n)$ and $z = \xi(z, n)$. If both f_1 and f_2 adhere to the requirements of CoL, their composition $z = f_2(f_1(x))$ also satisfy the setting of CoL, i.e., each z_i only depends on $x_{\leq i}$.

More important, CoL also supports compositionality, meaning that stacking multiple CoL layers also retains the characteristics of CoL. This property allows us to generalize the scope of CoL from the layer level to the model level.

2.3 Model

Definition 3 (Chain-of-Model) For a model θ with L layers, we define it as a *chain-of-model* (CoM) if each of its layers, from the first to the final, conforms to the chain-of-layer property.

Based on Definition 3, if a model satisfies the criteria of CoM, it also inherits all CoL properties, such as generality and causality. In other words, any model can be viewed as a kind of CoM (i.e., $n = 1$). And CoM can integrate multiple sub-models at different scales within one model, and enable us to scale up model from previous capabilities. This capability directly empowered foundation models with better extensibility and flexibility. In the next section, we will investigate how to build language model based on CoM framework.

3 Architecture

In this section, we describe the details about how to apply CoM into language model, including Linear, each module (e.g., embedding, self-attention, feedforward, normalization) within Transformer, and the objective function. Here, we term it as Chain-of-Language-Model (i.e., CoLM). Moreover, we further introduce a KV sharing mechanism based on CoLM framework, and term it as CoLM-Air, which provides better flexibility, such as prefilling speedup and seamless LLM switching, which always needs to re-calculate keys and values..

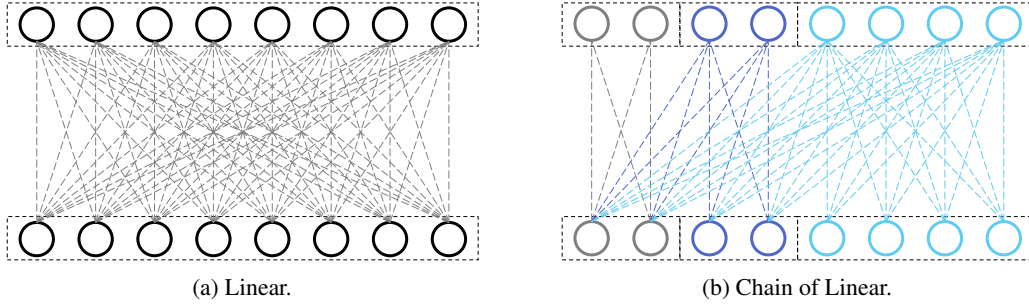


Figure 2: Comparisons between Linear and Chain-of-Linear layer.

3.1 Linear

Linear layer is a fundamental layer in neural network, which offers linear operation to transform the input features into the output features. The mathematic expression of linear layer is as $y = Wx + b$, where $W \in \mathbb{R}^{D_y \times D_x}$ and $b \in \mathbb{R}^{D_y}$. Here, each neuron of y is dense connected to the input feature x , so that we always need to compute all elements of W and b to preserve the integral semantics of y . Following the CoL setting, we introduce a new hyper-parameter into the linear layer, termed “Chain”, to determine the size of each chain in x and y as represented by CoR. Here, we set the hyper-parameter $\mathcal{C} = \{c_1, \dots, c_n\}$ as the base ratio to compute the dimension of each $x_i \in \mathbb{R}^{D_{x_i}}$ and $y_i \in \mathbb{R}^{D_{y_i}}$, where $D_{x_i} = \frac{c_i}{\sum \mathcal{C}} \cdot D_x$ and $D_{y_i} = \frac{c_i}{\sum \mathcal{C}} \cdot D_y$. Therefore, it is calculated as:

$$y = \parallel_{i=1}^n (y_i) = \parallel_{i=1}^n (W_i x_{\leq i} + b_i), \quad (1)$$

where $W_i \in \mathbb{R}^{D_{y_i} \times (\sum_{j=1}^i D_{x_j})}$, $b_i \in \mathbb{R}^{D_{y_i}}$ and \parallel denotes the concatenation operation. Given this design, the linear layer also satisfies the CoL criteria, and is identical to the standard linear when $n = 1$. Here, we term it as *Chain-of-Linear* or *Linear (Chain)* layer, which serves as the foundation to build language model. Figure 2 depicts the comparisons between Linear and Chain-of-Linear layer. To facilitate large-scale pre-training, we further provide an efficient implementation with a well-designed block-wise sparse kernel, optimizing both data access and all-reduce communication, as detailed in Appendix A.6.1.

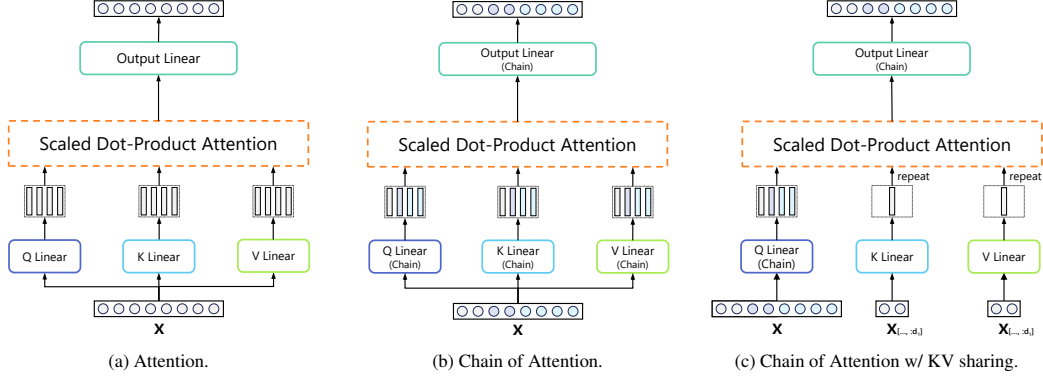


Figure 3: Differences between Attention, Chain of Attention and Chain of Attention with KV sharing.

3.2 Transformer

3.2.1 Multi-head Attention

In MHA, it applies linear operation to first obtain the query, key and value based on the input features, and split each into multiple heads. And then, we perform the attention function over each head (i.e., query, key and value) and concatenate all heads, followed by a linear transformation to produce the final output. Overall, the mathematic expression of MHA can be formulated as:

$$\text{MHA}(x) = O\left(\big\|_{i=1}^h \text{Attention}(q_i, k_i, v_i)\right) = O\left(\big\|_{i=1}^h \text{softmax}\left(\frac{q_i k_i^T}{\sqrt{d_k}}\right) v_i\right), \quad (2)$$

$$\text{where } \big\|_{i=1}^h q_i = Q(x), \big\|_{i=1}^h k_i = K(x), \big\|_{i=1}^h v_i = V(x),$$

where h is the head number, d_k is the head dimension, and $Q(\cdot)$, $K(\cdot)$, $V(\cdot)$, $O(\cdot)$ represent Linear layer to obtain query, key, value and output. To support CoR modeling in the MHA, we first directly replace all Linear layers (i.e., Q , K , V , O) by using Chain-of-Linear layers described in Sec 3.1. However, due to the dot product operations within $\text{Attention}(q, k, v)$, if a single head (i.e., q_i , k_i , v_i) contains information from 2 or more chains, its output will blend these information and then distribute to each chain, which cannot meet the setting of CoL. Therefore, we expect each head to only contain single scale information. To fulfill this, we design a simple trick that requires the sum of C is equal to h , so that each c_i represents how many heads will be allocated for chain i . This design enables each chain to retain their dedicated query, key, and value, ensuring our attention can satisfy the criteria of CoL, and we term it as *Chain-of-Attention* or *Attention (Chain)* module. Figure 3 illustrates the differences between Attention and Chain-of-Attention.

3.2.2 Feed-Forward Network

Feed-Forward Network (FFN) is another crucial component of the Transformer, which is composed of multiple linear layer interleaved with non-linear activation functions (e.g., ReLU [16] or SiLU [17]). To generalize the CoM framework into Transformer, FFN should also comply with the setting of CoL. To achieve this, we just need to substitute each Linear of FFN by using Chain-of-Linear layer, and adopt the same hyper-parameter C used as Attention (Chain) module. In this way, the output features of FFN can also conform to the requirements of CoR. The implementation and illustrative example about FFN can refer to Appendix A.6.3.

3.2.3 Normalization

Normalization is a technique to stabilize layer-wise training. In Transformer, we perform normalization [18, 19] at the feature level, before attention or FFN layers. To preserve the CoR properties, we apply the normalization function to each chain individually. Specifically, we use Root Mean Square Layer Normalization (RMSNorm) [19] for each chain to align with the LLaMA architecture. This chain-wise operation not only guarantees that nor-

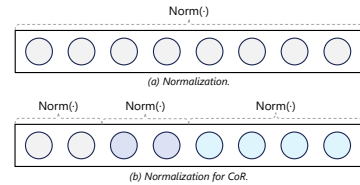


Figure 4: Normalization (CoR).

malized features maintain their causal structure but also slightly improves training efficiency by enabling group-wise operations, which can reduce the number of all-reduce communications. Figure 4 depicts the design of normalization layer.

3.2.4 Embedding

Embedding $\mathbf{E} \in \mathbb{R}^{V \times D}$ usually serves as the first layer in language model that converts the raw data \mathbf{I} as the input embeddings $x \in \mathbb{R}^D$. We do not involve any modifications at the Embedding during the training, but when encoding information at scale i , we just need to use embeddings $x[: \sum_{j=1}^i D_{x_j}]$ corresponding to the first i chains. A detailed example can be found in the Appendix A.6.4.

3.3 KV Sharing

Based on our above designs, Transformer can provide multiple capability across different models scales, enabling flexible deployment for various applications. However, in attention module, each chain also preserves its unique set of keys and values, and cannot bridge the connections between different scales. For example, when transitioning from a small language model (SLM) to a LLM for generation, it typically requires re-computing all keys and values for the preceding content. To this end, we further propose a ambitious technique, called *KV sharing*, where all keys and values are computed only in the first chain, and then shared to all other chains. If the number of keys and values is smaller than query heads, we follow the practices of GQA [20] by repeating keys and values to match the query heads, just as shown in Figure 3c. The design also meets the criteria of CoL¹. Although it slightly affect the performance, this design also exhibits some unique characteristics, e.g., prefilling speedup, and allowing us to seamlessly switch LMs at different scales from CoLM for generation. To the best of our knowledge, this is the first work to fulfill this functionality, and we term our CoLM with KV sharing as CoLM-Air. More detailed implementation are in Appendix A.6.2.

3.4 Objective Function

Based on the above designs of each component, the output feature $x \in \mathbb{R}^D$ of Transformer network will adhere to the CoR form $\xi(x, n)$. And finally, we apply a classification layer $\mathbf{W} \in \mathbb{R}^{D \times V}$, to map the feature dimension into the vocab size. Generally, we can adopt the cross entropy loss as the objective function for CoLM training. However, to enable multi-scale prediction, we also need to build individual classification heads for each scale. Motivated by MRL [21], we propose a multi-chain cross-entropy loss to compute each scale as: $\text{Loss}_i = \mathcal{L}(\mathbf{W}^i x_{\leq i})$, where $\mathbf{W}^i \in \mathbb{R}^{\sum_{j=1}^i (D_{x_j}) \times V}$ is a subset of \mathbf{W} , and \mathcal{L} is the cross entropy function. Although the computations of different logits (i.e., $\mathbf{W}^i x_{\leq i}$) can be shared without any additional cost, calculating multiple cross-entropy losses will still bring some computational overhead during the training. Due to resource limitations, we adopt a two-stage strategy inspired by MRL. We first pre-train the entire CoLM network using a standard cross-entropy loss calculated only on the output of the final, largest scale. This allows the model backbone to converge efficiently. After pre-training, we freeze the model backbone and efficiently fine-tune the classification heads for all scales simultaneously.

4 Experiments

4.1 Performance

In our experiments, we utilize the SlimPajama dataset [22] as the pre-training corpus, comprising 600 billion tokens. Here, we use LLama-2-tokenizer² to process corpus, with a vocab size of 32000. We leverage Fully Sharded Data Parallel (FSDP2) for distributed training, facilitated by the TorchTitan framework [23]. Our training is built upon a cluster of 32 NVIDIA A100 40GB GPUs with a gradient accumulation of 4, yielding an effective batch size of 1024. BFloat16 is adopted for training, and the sequence length is set as 4096 tokens. We choose AdamW optimizer [24] with a learning rate of 1.5×10^{-4} . Due to resource limitations, our model is pre-trained with 50K steps, nearly 200 billion tokens. For model configuration, we use chains \mathcal{C} and dimensions D to determine different models. The baseline use $\mathcal{C} = \{32\}$, and adopt same configuration as LLaMA-3.2-1B, so it is equivalent to

¹It can be considered as y_i is only conditioned on x_1 , which belongs to a subset of Definition 2.

²<https://huggingface.co/meta-llama/LLama-2-7b>

the standard language model. For CoLM-series, we use $\mathcal{C} = \{16, 16\}$ and $\mathcal{C} = \{8, 8, 8, 8\}$ ³. To align the baseline, we also increase the dimension of our CoLM as Chain-of-Linear takes fewer parameters than Linear layer. We use the EleutherAI Language Model Evaluation Harness [25] to evaluate models on commonsense tasks in a zero-shot setting. More details can refer to Appendix A.1.

Our results are reported in Table 1, and have these observations: 1) CoLM-Air with KV sharing mechanism will slightly affect model performance, as it only calculates keys and values within the first chain. However, this design introduces greater flexibility, which will be discussed in the following section; 2) using $\{16, 16\}$ can achieve better performance when compared with $\{8, 8, 8, 8\}$. When using $\{16, 16\}$ with a wider size, our method can achieve comparable performance to baseline. It can be considered as a trade-off as using more chains can offer more sub-models for usage. Overall, CoLM achieves results competitive with baseline while offering significantly faster prefilling speeds and greater flexibility, which will be discussed in subsequent sections.

Table 1: Performance of baseline and CoLM on commonsense reasoning tasks. Models are trained for 50K steps. Each task is reported by acc_norm metric.

Chain(\mathcal{C})	Dims	Params	HellaSwag \uparrow	Obqa \uparrow	WinoGranda \uparrow	ARC-e \uparrow	ARC-c \uparrow	Boolq \uparrow	Piqa \uparrow	Avg \uparrow
{32}	2048	1.10B	40.01	31.19	52.72	43.52	23.63	57.43	67.30	45.11
<i>CoLM</i>										
{16, 16}	2560	1.11B	40.25	31.39	52.41	43.73	23.81	58.01	67.30	45.27
{16, 16}	2048	0.86B	37.12	29.18	51.07	40.82	22.61	61.74	65.48	44.00
{8, 8, 8, 8}	3072	1.18B	38.63	31.99	51.62	42.80	23.89	56.70	66.00	44.51
{8, 8, 8, 8}	2048	0.74B	34.49	27.57	50.20	39.02	22.78	57.65	63.00	42.10
<i>CoLM-Air</i>										
{16, 16}	2560	1.11B	39.85	31.19	52.09	44.30	23.63	56.72	66.76	44.80
{16, 16}	2048	0.86B	36.82	28.77	51.62	49.19	22.70	61.31	65.94	43.90
{8, 8, 8, 8}	3072	1.18B	36.98	29.78	49.80	41.46	24.57	55.81	65.51	43.41
{8, 8, 8, 8}	2048	0.74B	33.77	27.97	49.41	38.93	22.61	57.65	62.10	41.77

4.2 Chain Expansion

Table 2: Performance of baseline versus expanded models on commonsense reasoning tasks.

Model	HellaSwag \uparrow	Obqa \uparrow	WinoGranda \uparrow	ARC-e \uparrow	ARC-c \uparrow	Boolq \uparrow	Piqa \uparrow	SciQ \uparrow	Avg \uparrow
<i>Tiny-LLaMA-v1.1</i>									
Baseline	61.47	36.62	59.43	55.47	32.68	55.99	73.56	84.20	57.43
+ Expansion	61.66	36.62	60.62	56.27	32.94	58.44	74.05	86.20	58.35
<i>LLaMA-3.2-1B</i>									
Baseline	63.78	36.42	59.83	60.56	36.03	63.70	74.37	88.40	60.39
+ Expansion	63.19	37.02	60.69	60.31	36.69	63.55	74.65	88.20	60.53

Given the generality and causality of CoM, any model can be regarded as a special case of CoM when chain is 1, and can be extended to multiple chains. Therefore, we introduce *Chain Expansion*, which use the well-trained model as the initial chain and then expand it with new additional chains. Conceptually, this is similar to progressive network architectures [12], that allowing us to increase additional capacity while retaining prior knowledge. However, CoLM has extended it to Transformer for training. Overall, *Chain Expansion* allows us to incrementally increase model scales based on the existing scales, avoid training from scratch.

To validate this point, we choose two LLaMA variants (i.e., TinyLLaMA-v1.1 [26] and LLaMA-3.2-1B [3]) as the first chain for expansion. Since these two models both use group query attention [20], where query number is 32 and kv head is 8. Following the Attention setting in Section 3.2.1, We set $c_1 = 32$ and subsequently introduce a second chain with $c_2 = 8$, incurring nearly 0.8B parameters. We perform zero initialization for the expanded classification head, so that its output logits is started from the original trained model. Due to resource limitations, our expanded models are trained with a batch size of 1024 for 20K steps, nearly 8 billion tokens. Besides, we also freeze the first chain to preserve the original knowledge and speedup training. More details can refer to Appendix A.1.

All of results are presented in Table 2. We achieve improvements of 0.92 and 0.14 over the Tiny-LLaMA-v1.1 and LLaMA-3.2-1B, respectively. Since LLaMa-3.2-1B is a stronger baseline, it

³Each chain is equal can slightly improve the training efficiency.

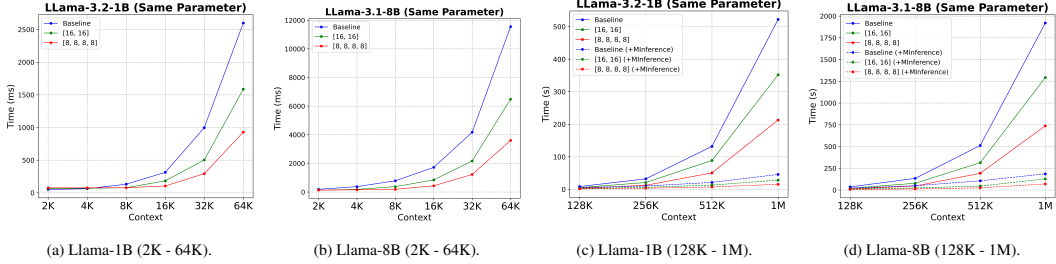


Figure 5: Comparisons on Prefilling speed on LLaMa-1B and LLaMa-8B settings (2K - 1M).

requires more computation to achieve significant gains but our method still can improve it under the limited computations. Overall, these results also indicate the effectiveness of our method in improving the baselines, even with constrained resources.

4.3 Elastic Inference

Elastic inference is to provide dynamic inference usage to meet the requirements from different deployment scenarios. Compared with previous methods [10, 11] that involved additional sampling policies, CoLM integrates multiple sub-models at varying scales within one model, enabling unified pre-training via a single forward pass. To obtain sub-models across different scales for prediction, we further pre-train our CoLM model using the loss function described in Sec 3.4. We choose the setting of CoLM-Air, where $\mathcal{C} = \{16, 16\}$ and Dims as 2048, for validation. Following prior work [21], we freeze the Transformer backbone and only train the classification head to accelerate the training. The results in Table 3 also highlight the potential of CoLM for enabling elastic inference.

Table 3: Elastic Inference to offer multiple sub-models by using different number of chains.

Scale	Params	HellaSwag	Obqa	WinoGranda	ARC-e	ARC-c	Boolq	Piqa	Avg
CoLM-Air, $\mathcal{C} = \{16, 16\}$, Dims = 2048									
Chain 1 + 2	0.86B	36.82	28.77	51.62	40.19	22.70	61.31	65.94	43.90
Chain 1	0.33B	29.31	25.75	52.96	34.50	22.01	62.23	61.15	41.13

4.4 Prefilling

Prefilling [27] is to generate keys and values in Transformer to understand the prompt of user instruction. Existing LLM paradigms always requires us to deploy the full model for calculation, introducing substantial computational costs. For CoLM-Air, all keys and values are computed within the first chain, so that we just need to calculate the first chain according to the causality of CoM. Please note the first chain is still a standard language model, we can further apply any inference-time techniques (e.g, MInference [28]). Here, we conduct two scales (LLama-3.2-1B and LLaMa-3.1-8B) for evaluation, with a context length from 2K to 1M. We choose two settings where $\mathcal{C} = \{8, 8, 8, 8\}$ and $\mathcal{C} = \{16, 16\}$ with similar parameters as baseline. Our results are shown in Figure 5. More experimental details and results can refer to Appendix A.2.1.

From Figure 5, We observe that CoLM-Air achieves faster prefilling speed when compared to LLaMa with similar parameter sizes. With the increasing of sequence length, CoLM-Air can achieve faster speedup at the prefilling. In particular, to process 1M tokens, CoLM-Air yields nearly $1.6\times$ and $3.0\times$ faster prefilling when \mathcal{C} is $\{16, 16\}$ and $\{8, 8, 8, 8\}$. And, when combining with MInference [28], CoLM-Air achieves at most $27\times$ speedup ($1920 \rightarrow 70.0$) in Figure 5d. These results all demonstrate that our CoLM-Air can effectively speedup prefilling.

4.5 Chain Tuning

Benefiting from the causality of CoM architecture, CoLM is composed of multiple chains, that each chain can leverage capability from previous chains. Therefore, based on this property, we propose *Chain Tuning*, which aims to fine-tune the latter chains while freezing the first few chains. Chain

Tuning can partially reduce the tuning costs and mitigate catastrophic forgetting via preserving the first few chains. Moreover, when we adopt the CoLM-Air setting and freeze the first chain, the keys and values from the fine-tuned models can seamlessly transfer to the original model, without any additional computations. Here, we utilize the expanded models (with two chains) trained in Section 4.2 as the baseline, and then fine-tune the expanded chain.

In addition, compared with Parameter-Efficient Fine-Tuning (PEFT) methods (e.g., LoRA [29]), our Chain Tuning is fully compatible with these PEFT methods, while further enable KV sharing to switch between the fine-tuned models and the original model. In other words, we can apply LoRA to the latter chains for even greater parameter efficiency. To demonstrate this capability, we choose GLUE benchmark [30] for a faster validation. We report the results within Table 4, and more experimental details can refer to Appendix A.1. From Table 4, we observe Chain-tuning can boost model performance by just fine-tuning approximately 42% of the model parameters.

Table 4: GLUE dev set results. Each task is reported by accuracy. “+ CT” means Chain Tuning, and the values in brackets mean the fine-tuned parameters and the total parameters.

Model	SST-2	COLA	MNLI	MRPC	QNLI	QQP	RTE	WNLI	Avg.
Baseline	67.89	32.98	35.87	49.02	51.77	51.65	54.87	43.66	48.46
+ LoRA	91.74	58.46	81.32	61.27	73.20	54.92	55.96	52.11	66.12
+ CT	92.32	54.87	82.26	62.01	82.87	55.22	59.21	53.52	67.79

4.6 Generative Performance

To assess performance on open-ended generation, we instruction-tuned our expanded TinyLLaMA-v1.1 model (from Section 4.2) on the Alpaca dataset [31] and evaluated it on MT-Bench [32]. The results in Table 5 show that our model, which leverages KV sharing, outperforms the similarly tuned baseline, particularly in reasoning and humanities tasks. This validates that the CoM framework and KV sharing mechanism are effective for enhancing instruction-following and conversational abilities.

Table 5: MT-Bench scores for instruction-tuned models. “Ours” refers to our expanded model. Both models are fine-tuned on the Alpaca dataset. Higher scores are better.

Model	Writing	Roleplay	Reasoning	Math	Coding	Extraction	STEM	Humanities	Avg.
TinyLLaMA-v1.1	2.85	3.20	2.05	1.10	1.30	1.40	2.35	3.20	2.18
Ours	2.75	3.40	2.70	1.00	1.10	1.30	2.35	3.65	2.28

5 Related Works

Foundation Models play a critical component in machine learning for effectively modeling relationships among data. From the perspective of architectural level, we posit that existing Large Foundation Models (LFM), can be conceptually decomposed into two parts, Base Architecture and Scaling Architecture. For Base Architecture, it refers to the backbone network (e.g., Transformer [7], Mamba [33], RWKV [34]) to perform feature transformation from the original data. For Scaling Architecture, it represents the scaling strategies (e.g., Increasing layers, dimensions or experts [35, 14]) to improve model capabilities. To some extent, CoM can also be regarded as a scaling architecture via increasing the number of chains, so that it can be applied to any base architectures. Moreover, compared with Dense or MoE architectures, the core idea of CoM is to integrate capabilities across varying scales and allow us to leverage inherited knowledge from previous models to improve scaling efficiency and avoid catastrophic forgetting.

Elastic Inference is to offer multiple sub-models tailored to handle diverse user instructions under different deployment constraints. Current LLM architectures can only provide a fixed inference capacity, and thus requires us to pre-train multiple models at different sizes, leading to inefficiencies. To this end, Once-for-All [10] build a super-net and then randomly sample a specialized sub-network for training. Some other works [11, 36] introduce some nested architecture designs into network blocks to offer elastic inference across different model granularity, but also require sampling strategy

of each sub-network. By incorporating causal dependencies among sub-models within CoR, CoM can integrate multi-scale training within a single forward pass, demonstrating better efficiency.

Continual Training is to continue the training process from the pre-trained models to boost its capability and avoid catastrophic forgetting. To leverage the capability of existing pre-trained models, the most simple method is using depth-up scaling method (i.e., increase the number of layers), such as Solar 10.7B [37]. Depth-up scaling can effectively reuse weights from previous layers for initialization, but fails to retain the output logits (i.e., probability) of the classification layers. CoM can be regarded as a width-up scaling method, which can perfectly preserves the knowledge of the entire model from the first layer to the final layer. Please note that CoM is also compatible with the depth-up scaling approach. Besides, Progressive neural network (PNN) [12] is a seminal work in continual training, which also motivates us to develop CoM. PNN extends the MLP layers by incorporating new parameters via lateral connections, and then apply to the reinforcement learning tasks. Compared to PNN, CoM extends its scope to a broader range and effectively demonstrates the feasibility of using Transformers for pre-training.

6 Conclusion

In this paper, we introduce the concept of “*Chain-of-Representation*” (CoR) to encode multi-scale information within hidden representation. Subsequently, we further propose “*Chain-of-Model*” (CoM) learning to capture the causal relationship across varying scales among CoR features of each layer within the model. The model built upon the CoM framework can provide multi-scale functionality and retain capability from previous scales for usage. Therefore, we devise “*Chain-of-Language-Model*” (CoLM) by integrating the Chain-of-Model (CoM) paradigm into each layer of the Transformer architecture. As a result, CoLM can offer multi-scale sub-models at varying scales within a unified language model. Moreover, we propose a KV sharing mechanism that computes all keys and values within the first chain. This design significantly unleashes the flexibility and extensibility of language models, including prefilling acceleration, dynamic LLM switching and so on. Experimental results across multiple datasets demonstrate that our CoLM family matches the performance of standard Transformers, while offering superior extensibility and adaptability across a wide range of scenarios, paving a new way to build next-generation foundation models.

7 Limitations

While Chain-of-Model (CoM) introduces a flexible and efficient paradigm for scaling language models, some challenges and areas for future exploration remain. First, the Chain-of-Linear layer, in its current form, presents complexities for Tensor Parallelism (TP), a key technique for training extremely large models. Although compatible with other parallelism strategies, optimizing CoM for TP-centric infrastructures requires further investigation. Second, in the CoLM-Air architecture, the first chain’s quality is critical, as it generates the keys and values for all subsequent chains. Determining the optimal configuration and training strategy for this foundational chain is crucial for maximizing the overall model’s potential. We provide a more detailed discussion of these limitations and outline potential future directions in Appendix A.5.

References

- [1] OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023.
- [2] Google Gemini Team. Gemini: A family of highly capable multimodal models. *CoRR*, abs/2312.11805, 2023.
- [3] Meta Llama Team. The llama 3 herd of models. *CoRR*, abs/2407.21783, 2024.
- [4] Alibaba Qwen Team. Qwen2.5 technical report. *CoRR*, abs/2412.15115, 2024.
- [5] Deepseek Team. Deepseek-v3 technical report. *CoRR*, abs/2412.19437, 2024.
- [6] Mistral AI Team. Mixtral of experts. *CoRR*, abs/2401.04088, 2024.

- [7] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008, 2017.
- [8] 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. *CoRR*, abs/2001.08361, 2020.
- [9] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training compute-optimal large language models. *CoRR*, abs/2203.15556, 2022.
- [10] Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. Once-for-all: Train one network and specialize it for efficient deployment. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.
- [11] Devvrit, Sneha Kudugunta, Aditya Kusupati, Tim Dettmers, Kaifeng Chen, Inderjit S. Dhillon, Yulia Tsvetkov, Hannaneh Hajishirzi, Sham M. Kakade, Ali Farhadi, and Prateek Jain. Mat-former: Nested transformer for elastic inference. *CoRR*, abs/2310.07707, 2023.
- [12] Andrei A. Rusu, Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *CoRR*, abs/1606.04671, 2016.
- [13] Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual learning: Theory, method and application. *IEEE Trans. Pattern Anal. Mach. Intell.*, 46(8):5362–5383, 2024.
- [14] Robert A. Jacobs, Michael I. Jordan, Steven J. Nowlan, and Geoffrey E. Hinton. Adaptive mixtures of local experts. *Neural Comput.*, 3(1):79–87, 1991.
- [15] James Kirkpatrick, Razvan Pascanu, Neil C. Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural networks. *CoRR*, abs/1612.00796, 2016.
- [16] Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10), June 21-24, 2010, Haifa, Israel*, pages 807–814, 2010.
- [17] Noam Shazeer. GLU variants improve transformer. *CoRR*, abs/2002.05202, 2020.
- [18] Lei Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *CoRR*, abs/1607.06450, 2016.
- [19] Biao Zhang and Rico Sennrich. Root mean square layer normalization. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 12360–12371, 2019.
- [20] Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. GQA: training generalized multi-query transformer models from multi-head checkpoints. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 4895–4901, 2023.
- [21] Aditya Kusupati, Gantavya Bhatt, Aniket Rege, Matthew Wallingford, Aditya Sinha, Vivek Ramanujan, William Howard-Snyder, Kaifeng Chen, Sham M. Kakade, Prateek Jain, and Ali Farhadi. Matryoshka representation learning. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.

- [22] Daria Soboleva, Faisal Al-Khateeb, Robert Myers, Jacob R Steeves, Joel Hestness, and Nolan Dey. SlimPajama: A 627B token cleaned and deduplicated version of RedPajama. <https://cerebras.ai/blog/slimpajama-a-627b-token-cleaned-and-deduplicated-version-of-redpajama>, June 2023.
- [23] Wanchao Liang, Tianyu Liu, Less Wright, Will Constable, Andrew Gu, Chien-Chin Huang, Iris Zhang, Wei Feng, Howard Huang, Junjie Wang, Sanket Purandare, Gokul Nadathur, and Stratos Idreos. TorchTitan: One-stop pytorch native solution for production ready LLM pre-training. *CoRR*, abs/2410.06511, 2024.
- [24] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- [25] 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. The language model evaluation harness, 07 2024.
- [26] Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small language model. *CoRR*, abs/2401.02385, 2024.
- [27] Amey Agrawal, Ashish Panwar, Jayashree Mohan, Nipun Kwatra, Bhargav S. Gulavani, and Ramachandran Ramjee. SARATHI: efficient LLM inference by piggybacking decodes with chunked prefills. *CoRR*, abs/2308.16369, 2023.
- [28] Huiqiang Jiang, Yucheng Li, Chengruidong Zhang, Qianhui Wu, Xufang Luo, Surin Ahn, Zhenhua Han, Amir Abdi, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. Minference 1.0: Accelerating pre-filling for long-context llms via dynamic sparse attention. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024.
- [29] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022.
- [30] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.
- [31] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.
- [32] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- [33] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *CoRR*, abs/2312.00752, 2023.
- [34] Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Leon Derczynski, Xingjian Du, Matteo Grella, Kranthi Kiran GV, Xuzheng He, Haowen Hou, Przemysław Kazienko, Jan Kocon, Jiaming Kong, Bartłomiej Koptyra, Hayden Lau, Jiaju Lin, Krishna Sri Ipsit Mantri, Ferdinand Mom,

- Atsushi Saito, Guangyu Song, Xiangru Tang, Johan S. Wind, Stanislaw Wozniak, Zhenyuan Zhang, Qinghua Zhou, Jian Zhu, and Rui-Jie Zhu. RWKV: reinventing rnns for the transformer era. In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 14048–14077, 2023.
- [35] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarczyk, Andy Davis, Quoc V. Le, Geoffrey E. Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *CoRR*, abs/1701.06538, 2017.
- [36] Abhinav Shukla, Sai Vemprala, Aditya Kusupati, and Ashish Kapoor. Matmamba: A matryoshka state space model. *CoRR*, abs/2410.06718, 2024.
- [37] Sanghoon Kim, Dahyun Kim, Chanjun Park, Wonsung Lee, Wonho Song, Yunsu Kim, Hyeonwoo Kim, Yungi Kim, Hyeonju Lee, Jihoo Kim, Changbae Ahn, Seonghoon Yang, Sukyung Lee, Hyunbyung Park, Gyoungjin Gim, Mikyoung Cha, Hwalsuk Lee, and Sunghun Kim. SOLAR 10.7b: Scaling large language models with simple yet effective depth up-scaling. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Industry Track, NAACL 2024, Mexico City, Mexico, June 16-21, 2024*, pages 23–35. Association for Computational Linguistics, 2024.
- [38] Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- [39] Yinmin Zhong, Shengyu Liu, Junda Chen, Jianbo Hu, Yibo Zhu, Xuanzhe Liu, Xin Jin, and Hao Zhang. Distserve: Disaggregating prefill and decoding for goodput-optimized large language model serving. In *18th USENIX Symposium on Operating Systems Design and Implementation, OSDI 2024, Santa Clara, CA, USA, July 10-12, 2024*, pages 193–210. USENIX Association, 2024.
- [40] Pratyush Patel, Esha Choukse, Chaojie Zhang, Aashaka Shah, Íñigo Goiri, Saeed Maleki, and Ricardo Bianchini. Splitwise: Efficient generative LLM inference using phase splitting. In *51st ACM/IEEE Annual International Symposium on Computer Architecture, ISCA 2024, Buenos Aires, Argentina, June 29 - July 3, 2024*, pages 118–132. IEEE, 2024.
- [41] Ruoyu Qin, Zheming Li, Weiran He, Mingxing Zhang, Yongwei Wu, Weimin Zheng, and Xinran Xu. Mooncake: A kvcache-centric disaggregated architecture for LLM serving. *CoRR*, abs/2407.00079, 2024.
- [42] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott E. Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 1–9. IEEE Computer Society, 2015.
- [43] William Brandon, Mayank Mishra, Aniruddha Nrusimha, Rameswar Panda, and Jonathan Ragan-Kelley. Reducing transformer key-value cache size with cross-layer attention. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024.
- [44] Zayd Muhammad Kawakibi Zuhri, Muhammad Farid Adilazuarda, Ayu Purwarianti, and Alham Fikri Aji. MLKV: multi-layer key-value heads for memory efficient transformer decoding. *CoRR*, abs/2406.09297, 2024.
- [45] Yutao Sun, Li Dong, Yi Zhu, Shaohan Huang, Wenhui Wang, Shuming Ma, Quanlu Zhang, Jianyong Wang, and Furu Wei. You only cache once: Decoder-decoder architectures for language models. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024.
- [46] David Raposo, Samuel Ritter, Blake A. Richards, Timothy P. Lillicrap, Peter Conway Humphreys, and Adam Santoro. Mixture-of-depths: Dynamically allocating compute in transformer-based language models. *CoRR*, abs/2404.02258, 2024.

- [47] Dmitry Lepikhin, HyukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. Gshard: Scaling giant models with conditional computation and automatic sharding. *CoRR*, abs/2006.16668, 2020.
- [48] Damai Dai, Chengqi Deng, Chenggang Zhao, R. X. Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng, Xingkai Yu, Y. Wu, Zhenda Xie, Y. K. Li, Panpan Huang, Fuli Luo, Chong Ruan, Zhifang Sui, and Wenfeng Liang. Deepseekmoe: Towards ultimate expert specialization in mixture-of-experts language models. *CoRR*, abs/2401.06066, 2024.

A Appendix

A.1 Experimental Details

Datasets SlimPajama [22] is a deduplicated, multi-corpora, open-source dataset for training LLMs. We use SlimPajama dataset exclude books corpus as the pre-training corpus, result in approximately 600 billion tokens. The composition of our used dataset is shown in Table 6.

Dataset	Ratio
Commoncrawl	52.2%
C4	26.7%
GitHub	5.2%
ArXiv	4.6%
Wikipedia	3.8%
StackExchange	3.3%

Table 6: Dataset statistical of SlimPajama dataset.

Hyperparameters Table 7 presents the hyperparameters that were used in pretraining. Besides, FlashAttention-2 [38] is applied to speed up attention computation. For the Expansion experiment, we set the training step to 20,000, while preserving the same configuration as the pre-training. Table 8 presents the hyperparameters used for Chain Tuning.

Evaluation In this project, we utilized the Eleuther AI Language Model Evaluation Harness library [25]. We used the default evaluation configuration. For chain tuning, we set the prompt template to the default from the evaluation scripts.

Table 7: Hyperparameters used in pretraining.

Hyperparameter	Value
Learning Rate	$2e - 4$
Sequence Length	4096
Batch Size Per Device	8
Gradient Accumulation Steps	8
Optimizer	Adamw
LR Scheduler Type	Linear
Warmup Steps	2000
Max Norm	1.0
Activation checkpointing	full
Mixed Precision for Parameters	bfloat16
Mixed Precision for Reductions	float32

Table 8: Hyperparameters used in Chain Tuning.

Hyperparameter	Value
Learning Rate	$2e - 4$
Epochs	1
Batch Size Per Device	4
Gradient Accumulation Steps	32
Optimizer	Adamw
LR Scheduler Type	Cosine
Warmup Ratio	0.1
Sequence Length	4096

Model Configuration Table 9 presents the detailed configurations for all the models used in our experiments. Each model is identified by its name and characterized by configurations such as the Chain setting, dimension (dim), hidden dimension (hidden dim), the number of layers, attention

Table 9: Model Configuration.

Model Name	Chains	# Params	Dim	Hidden Dim	# Layers	# Heads	# KV Heads
expansion-llama3.2-1B	{32, 8}	1.93B	2560	10240	16	40	10
expansion-tinyllama	{32, 8}	1.44B	2560	7040	22	40	5
Llama	{32}	1.10B	2048	8192	16	32	8
CoLM_16-16	{16, 16}	0.86B	2048	8192	16	32	8
CoLM_8-24	{8, 24}	0.92B	2048	8192	16	32	8
CoLM_8-8-16	{8, 8, 16}	0.80B	2048	8192	16	32	8
CoLM_8-8-8-8	{8, 8, 8, 8}	0.74B	2048	8192	16	32	8
CoLM_16-16-same-params	{16, 16}	1.11B	2560	8192	16	32	8
CoLM_8-8-8-8-same-params	{8, 8, 8, 8}	1.18B	3072	8192	16	32	8

heads, key-value heads and the number of parameters. Figure 12 illustrates the structure of the Chain of Transformer, designed for multi-scale representation learning.

A.2 Full Results

A.2.1 Prefilling

Our prefilling is evaluated on a single NVIDIA A100 GPU with a batch size of 1. Specifically, from 128K to 1M, we perform chunk-wise prefilling (e.g., chunk operation for layer calculation) to avoid the out-of-memory (OOM) issue, and we also evaluate and validate our method when combining with MInference [28]. Here, we evaluate six configurations by using different chain settings ({16, 16}, {8, 8, 8, 8}, {8, 8, 16}, {8, 24}) and increase the number of layers to ensure the parameter sizes are comparable. Table 10 and Table 11 report all results of our CoLM-Air setting (1B and 8B) of a sequence length from 2K to 64K. And Table 12 and Table 13 report all results of our CoLM-Air setting (1B and 8B) of a sequence length from 128K to 1M. From these results, we can find that CoLM can significantly speed up the prefilling stage in calculating long-context instructions.

Moreover, our CoLM is also extremely suitable for PD disaggregation architecture [39–41], as it only requires us to deploy the entire weights of the first chain on the prefilling server to calculate KV caches, which significantly reduces the computation loads of the prefilling server.

Table 10: Prefilling results on LLama-3.2-1B setting (2K - 64K). All results are reported in milliseconds. D_{x_1} represents the dimension of the first chain.

Method	L	$D(D_{x_1})$	Length					
			2K	4K	8K	16K	32K	64K
LLama-3.2-1B	16	2048 (2048)	51.0	62.0	132.0	313.0	997.0	2,600
<i>Same configuration (CoLM-Air)</i>								
{16,16}	16	2048 (1024)	47.0	47.6	60.0	140.0	380.0	1,189
{8, 8, 8, 8}	16	2048 (512)	48.0	50.0	51.0	71.0	193.0	613.0
<i>Similar parameters (CoLM-Air)</i>								
{16,16}	22	2048 (1024)	66.0	69.0	78.0	183.0	503.0	1,587
{8, 8, 8, 8}	26	2048 (512)	76.0	76.0	78.0	105.0	293.0	929.0
{8,8,16}	24	2048 (512)	68.0	69.0	73.0	99.0	272.0	863.0
{8,24}	21	2048 (512)	62.0	60.0	65.0	87.0	243.0	768.0

A.2.2 CIFAR-10 MLP Experiments

We compare a baseline fully-connected MLP against several Chain-of-Model MLP (CoM-MLP) variants on CIFAR-10. The baseline model uses a 4-layer architecture with hidden sizes [768, 1024, 512, 256] (an input embedding of dimension 768, followed by three dense layers and a 10-classes output). CoM-MLP models augment each layer with multiple sequential conditional modules specified by the “Chain-setting” vector. Models were trained using the Adam optimizer with an initial

Table 11: Prefilling results on LLama-3.1-8B setting (2K - 64K). All results are reported in seconds. D_{x_1} represents the dimension of the first chain. D_{x_1} represents the dimension of the first chain.

Method	L	$D(D_{x_1})$	Length					
			2K	4K	8K	16K	32K	64K
LLama-3.1-8B	32	4096 ₍₄₀₉₆₎	0.19	0.38	0.79	1.72	4.18	11.55
<i>Same configuration (CoLM-Air)</i>								
{16,16}	32	4096 ₍₂₀₄₈₎	0.09	0.14	0.28	0.64	1.64	4.89
{8, 8, 8, 8}	32	4096 ₍₁₀₂₄₎	0.09	0.10	0.12	0.28	0.72	2.30
<i>Same parameters (CoLM-Air)</i>								
{16,16}	43	4096 ₍₂₀₄₈₎	0.13	0.18	0.37	0.85	2.17	6.49
{8, 8, 8, 8}	52	4096 ₍₁₀₂₄₎	0.13	0.16	0.19	0.43	1.06	3.61
{8,8,16}	48	4096 ₍₁₀₂₄₎	0.13	0.15	0.17	0.34	0.94	3.35
{8,24}	40	4096 ₍₁₀₂₄₎	0.11	0.14	0.15	0.85	2.17	2.82

Table 12: Prefilling results on LLama-3.2-1B setting (128K - 1M). All results are reported in seconds. The value in the bracket is using the MInference trick. D_{x_1} represents the dimension of the first chain.

Method	L	$D(D_{x_1})$	Length			
			128K	256K	512K	1M
LLama-3.2-1B	16	2048 ₍₂₀₄₈₎	8.6 (5.5)	32.4 (10.3)	132.0 (21.4)	522.0 (45.8)
<i>Same configuration (CoLM-Air)</i>						
{16, 16}	16	2048 ₍₁₀₂₄₎	4.0 (2.1)	15.4 (4.5)	63.2 (9.7)	265.0 (20.1)
{8, 8, 8, 8}	16	2048 ₍₅₁₂₎	1.9 (1.0)	7.5 (2.2)	30.6 (4.8)	124.0 (9.9)
<i>Same parameters (CoLM-Air)</i>						
{16, 16}	22	2048 ₍₁₀₂₄₎	5.4 (2.9)	21.0 (6.2)	88.8 (13.1)	352.0 (28.4)
{8, 8, 8, 8}	26	2048 ₍₅₁₂₎	3.2 (1.6)	12.5 (3.5)	50.7 (7.4)	213.0 (15.5)
{8, 8, 16}	24	2048 ₍₅₁₂₎	2.9 (1.5)	11.3 (3.1)	46.6 (6.7)	188.0 (14.5)
{8, 24}	21	2048 ₍₅₁₂₎	2.6 (1.3)	9.9 (2.8)	40.7 (6.0)	165.0 (13.3)

learning rate of 10^{-3} , a batch size of 64, for a total of 20 epochs. Reported accuracies are the mean \pm standard deviation over 3 independent runs.

Table 14 summarizes the test accuracy and parameter count for the baseline MLP and CoM-MLP configurations. While all models achieve similar accuracy (around 54–55%), the CoM-MLP with Chain-setting {4, 12} achieves the best accuracy ($55.11 \pm 0.21\%$) with fewer parameters compared to the dense baseline.

A.3 Training Loss Curves

We conducted several small-scale experiments (batch size 8, 15K steps) to study how different CoLM configurations and settings affect training convergence. The goal was to understand the impact of CoM design choices on convergence speed and final loss. Figures 6 plot the training loss over time for each experiment.

Comparison of CoLM Configurations Figure 6(a) shows the training loss curves for four chain configurations: equal four-chain {8, 8, 8, 8}, two 16-unit chains {16, 16}, and two mixed settings {8, 24} and {8, 8, 16}. The 16-16 configuration converges more quickly and reaches the lowest loss, indicating more efficient training, which appears to yield faster convergence and lower loss under these conditions, suggesting a trade-off in chain design.

Table 13: Prefilling results on LLama-3.1-8B setting (128K - 1M). All results are reported in seconds. The value in the bracket is using the MInference trick. D_{x_1} represents the dimension of the first chain.

Method	L	$D(D_{x_1})$	Length			
			128K	256K	512K	1M
LLama-3.1-8B	32	4096 (4096)	36.0 (21.6)	134.0 (49.5)	512.0 (106.0)	1920 (186.0)
<i>Same configuration (CoLM-Air)</i>						
{16, 16}	32	4096 (2048)	15.6 (6.4)	60.0 (13.3)	244.8 (28.5)	927.0 (89.0)
{8, 8, 8, 8}	32	4096 (1024)	7.3 (2.5)	28.1 (5.1)	117.3 (10.9)	472.0 (44.7)
<i>Same parameters (CoLM-Air)</i>						
{16, 16}	43	4096 (2048)	20.6 (10.8)	77.7 (23.5)	314.8 (45.6)	1293 (128.0)
{8, 8, 8, 8}	52	4096 (1024)	12.0 (6.2)	46.4 (12.5)	192.4 (25.4)	738 (70.1)
{8, 8, 16}	48	4096 (1024)	11.0 (5.9)	40.5 (10.7)	167.4 (23.9)	680 (65.0)
{8, 24}	40	4096 (1024)	9.2 (4.2)	35.4 (9.7)	146.0 (19.6)	591 (56.1)

Table 14: CIFAR-10 classification accuracy (mean \pm std) and parameter count for the baseline MLP and CoM-MLP variants.

Model	Chains	#Params	Accuracy (%)
MLP (Dense)	-	3,806,218	54.93 \pm 0.23
CoM-MLP	{4, 4, 4, 4}	3,263,754	54.60 \pm 0.16
CoM-MLP	{8, 8}	3,443,978	54.68 \pm 0.04
CoM-MLP	{4, 12}	3,534,090	55.11 \pm 0.21
CoM-MLP	{4, 4, 8}	3,353,866	54.55 \pm 0.36

Loss Across Different Heads We also examined head-wise training by optimizing each chain output head together in the {8, 8, 8, 8} CoLM setting. Figure 6(c) plots the training loss for each of the four heads. The results reveal that larger chains effectively leverage the outputs of preceding chains, resulting in greater capacity and improved performance. This trend demonstrates that deeper chains benefit from the cumulative computations of earlier chains, enhancing their ability to model complex features. The intermediate heads show intermediate performance, reflecting a gradual increase in capacity along the chain depth. These findings highlight the hierarchical structure of CoLM, where later chains achieve stronger capacity by building on the representations provided by earlier ones.

Effect of KV Sharing Figure 6(c) compares the {8, 8, 8, 8} setting with and without the KV-sharing mechanism (CoLM-Plus). While KV sharing results in a slightly higher loss for the same head, it

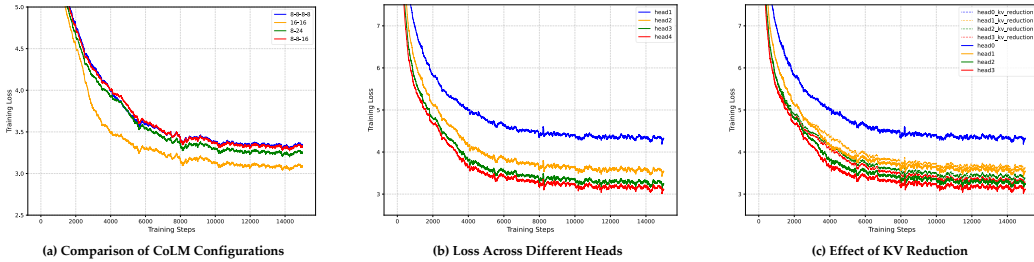


Figure 6: (a) Training loss for different CoLM chain configurations. (b) Training loss for each output head in the {8, 8, 8, 8} model (heads 1–4). (c) Training loss for the {8, 8, 8, 8} configuration, with versus without KV sharing.

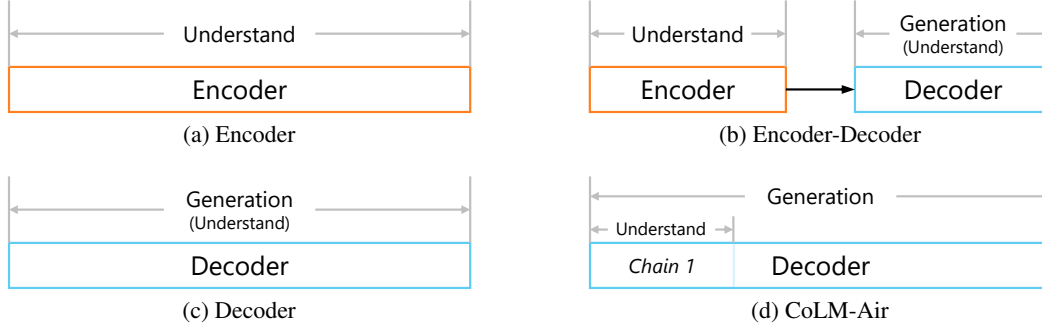


Figure 7: Illustrative comparisons between Encoder, Encoder-Decoder, Decoder and CoLM-Air, in term of understanding and generation capabilities.

significantly improves prefilling speed, as discussed in Section 4.4. This trade-off highlights the efficiency gains of KV sharing without a substantial impact on performance.

A.4 Discussion

Question 1 *Chain-of-model architecture seems to be a major improvement in the width of language models (i.e., dimension). Can it be applied to the depth of language models (i.e., layers)?*

Yes, the existing mechanism of chain-of-model design requires us to first determine the depth of large foundation models, and then allows us to extend its capability from the perspective of the dimension (i.e., width). Actually, we also conduct preliminary experiments to explore the feasibility of applying the idea of chain-of-model across the depth (i.e., layers) of the language model. Here, when applying CoM into the model depth, we also require it to meet two characteristics to ensure the alignment with the settings of our current version (i.e., CoLM):

1. The model should encompass multiple sub-models to enable dynamic selection and inference;
2. Each sub-model can seamlessly reuse past states (e.g., keys and values in Transformer) that are generated by different sub-models.

To meet these requirements, we have made two major designs based on the Transformer architecture: (a) we select some intermediate layers and then append the classifier to each as the sub-models; (b) we limit keys and values are only produced by the first sub-model (the first few bottom layers), which are shared across the subsequent layers. In our internal experiments, the setting (a) can guarantee that the model can be optimized and reach convergence, while we can also obtain multiple nested sub-models. Previous researches (e.g., GoogleNet [42]) also validate the availability of such a solution. But this design also introduces additional complexity, as we need to maintain multiple classifiers, which is extremely cumbersome and redundant. Although we can tie the weights of each classifier together to reduce model parameters, it will also cause interference to the model performance and cannot save any computations to calculate each classification head. However, when further integrated with the setting (b), which delivers the keys and values from the bottom layers to the top layers, we find that it severely degrades model performance. In other words, only using bottom-level states (i.e., keys and values) for all Transformer layers is not enough to create high-level semantics for next-token prediction. Existing attempts [43, 44] also indicate that we can only restrict key and value sharing to adjacent layers, which identifies the importance of depth to abstract high-level feature states (i.e., keys and values). Besides, some recent works (e.g., YOCO [45]) have explored reusing features from the outputs of half of the layers, which introduces modifications that combine efficient self-attention and cross-attention mechanisms. However, this technique is restricted to KV cache reduction during prefilling and does not enable multiple sub-model inference. Some other works (e.g., MoD [46]) introduce routers to dynamically skip certain layers for faster inference. But this design requires us to reset the keys and values to zero for the skipped layers, which potentially introduces discrepancies between training and inference. In our CoM method, we can retain all keys and values at each layer, enabling inference at multiple scales. Generally speaking, to abstract high-level features in neural network, depth usually matters more than width and preserving invariant depth (i.e., layers) for

foundation models is extremely necessary. This can also explain why our current version can only be applied to width, rather than depth. We will continue exploring solutions that extend the idea of chain-of-model to both width and depth in the future.

Question 2 *What are the differences between the Chain-of-Model (CoM) and Mixture-of-Experts (MoE) architectures?*

Table 15: Comparisons between MoE and CoM from different aspects.

	Mixture-of-Expert (MoE)	Chain-of-Model (CoM)
Expert Capability	Equivalent	Weak \rightarrow Strong
Expert Range	FFN	Full Model
Expert Activation	Sparse	Nested

The concept of Mixture of Experts (MoE) was first proposed by [14], which aims to create multiple expert networks and then use a gating mechanism for selective activation. As Transformer has exhibited huge potential in scalability, some works [35, 47] have investigated how to apply the MoE architecture into the FFN layers within Transformer to build the giant network. On the basis of this, some works [6, 48] have successfully scaled up language models with MoE architectures to a substantial size. We deem that MoE can also be considered as a kind of scaling architecture. Existing MoE architectures in Transformer are designed to create multiple experts with similar capabilities at the FFN layers, while only a few experts are activated for each token prediction. In the CoM structure, it aims to create a sequence of nested experts, progressing from weak to strong capability. Specifically, each expert contains the capability of all its preceding experts (i.e., weaker experts). In Table 15, we also illustrate the comparisons between CoM and MoE architectures. Therefore, CoM and MoE are designed from different perspectives to construct experts. In other words, CoM is completely orthogonal to MoE, meaning both of them can be applied within the same architecture and inherit advantages from each. We will explore how to combine CoM and MoE architecture, and leave it as future work.

Question 3 *When do we need the Chain-of-Model architecture?*

Generally, Chain-of-Model (CoM) architecture can be considered as an innovative design to optimize the scaling architecture (e.g., increasing model capabilities) of foundation models, rather than altering the backbone network (e.g., Transformer). Therefore, it can be generalized to any model with different architectures, such as CNNs or parallelized RNNs. However, this does not mean we can apply CoM architecture to any scenario. Based on the designs of CoM architecture, it can be considered as a composition of multiple nested sub-models. Each sub-model will contribute to the capability of the subsequent models. More specifically, the first chain (sub-model) plays a critical role in CoM architectures as it directly determines the understanding capability of the entire model. Therefore, if the capacity of the model is not sufficient to exhibit the power of the scaling law or lacks enough generalization, it is unnecessary to use the architecture of the chain-of-model method. In other words, the CoM structure is better suitable for large-scale model sizes, rather than smaller ones. This also indicates that CoM is an advanced design to optimize the scaling architecture.

A.5 Limitations

From the perspective of methodology, Chain-of-Model (CoM) introduces an innovative and brand-new solution to revolutionize the scaling architecture of large-scale foundation models. However, CoM still remains some ongoing challenges when applying it to practical applications:

- Although the CoM architecture offers exceptional flexibility in utilizing LLMs, it also poses new challenges for training a large-scale CoM architecture at the infrastructure level. Specifically, Chain-of-Linear layer is compatible with data, pipeline, and context parallelism (DP, PP, CP), but not well-suited for tensor parallelism (TP). In the naive implementation (Please see Figure 2b or Algorithm 1), the Chain-of-Linear layer can be viewed as a composition of multiple sub-linear layers. As a result, if we apply tensor parallelism to each sub-linear within the Chain-of-Linear

layer, it will significantly increase the number of all-reduce operations and data access when compared with the standard Linear. Although block-wise sparse kernel optimization can help us to alleviate these issues, we still need to devise a new pipeline mechanism when handling large tensor inputs. To some degree, Chain-of-Linear layer can be considered as an imbalanced split-wise tensor, which requires us to specifically design a tensor parallelism mechanism for it. Here, we leave it as future work.

- Empowered by chain-of-model design, CoLM is a novel architecture that combines multiple nested language models at varying scales. Benefiting from such a design, each model will inherit all its nested sub-models and extend their capability. In other words, each nested model will affect its subsequent model. More specifically, CoLM further introduces a KV-reduction mechanism to guarantee that the understanding parts (e.g., keys and values in Transformer) are all produced by the first chain. That means the first chain plays a very important role in building CoLM as it determines the understanding capability of the entire model. So, compared to the other chains, we expect the first chain to pay more attention to expressive and generalized content states (i.e., keys and values), in relative to its generation. Therefore, how to explore the optimal setting for the first chain could be very important in the CoLM framework. However, due to resource limitations, we leave it as future works.

Algorithm 1 Pseudo Code for Chain-of-Linear Layer (Linear_chain.py)

```
class Linear(nn.Module):

    def __init__(
        self, dim: int, hidden_dim: int, chains: List[int], bias = False
    ):
        """
        dim:          Dimension of input features
        hidden_dim:   Dimension of output features
        chains:       Base of each chain
        bias:         whether using bias or not
        """
        super().__init__()
        self.num_chain = len(chains)
        self.in_dims = [dim * c // sum(chains) for c in np.cumsum(chains)]
        self.out_dims = [hidden_dim * c // sum(chains) for c in chains]
        self.mlps = nn.ModuleList([
            nn.Linear(indim, outdim, bias=bias)
            for indim, outdim in zip(self.in_dims, self.out_dims)
        ])

    def forward(self, x: torch.Tensor, chain: int = None) -> torch.Tensor:
        """
        x:          Input features, shape = (B, L, D)
        chain:      The number of activated chains ([0, num_chain - 1])
        """
        if chain is None:
            chain = self.num_chain - 1

        assert (
            0 <= chain < self.num_chain
        ), f"The chain id should be in [0, {self.num_chain})."
        outputs = [
            self.mlps[i](x[..., :indim])
            for i, indim in enumerate(self.in_dims[:chain + 1])
        ]
        return torch.cat(outputs, dim=-1)
```

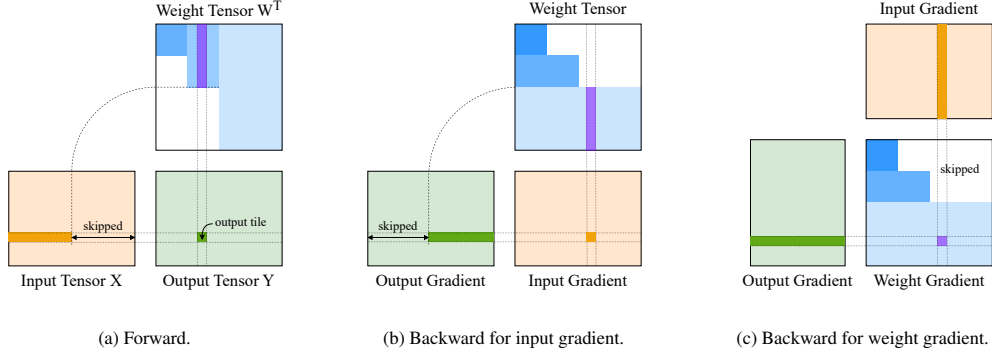


Figure 8: Block-sparse GeMM kernel for Linear (Chain) layer.

Table 16: Speed comparisons between different kernels. The “fwd”, “bwd_w” and “bwd_i” represent the speed of forward, backward (weight) and backward (inputs). Chain-of-Linear (Naive) refers to Algorithm 1, and Chain-of-Linear (Triton) corresponds to our block-wise sparse kernel implementation.

Implementation	Dim	Params	Computation		
			fwd	bwd _w	bwd _i
$\mathcal{C} = \{8, 8, 8, 8\}$, <i>tensor size is (4096, 4096).</i>					
Linear	4096	167K	0.67	0.62	0.64
Chain-of-Linear (Naive)	4096	104K	0.58	0.49	0.72
Chain-of-Linear (Triton)	4096	104K	0.47	0.46	0.56
$\mathcal{C} = \{4, 4, 4, 4, 4, 4, 4, 4\}$, <i>tensor size is (8192, 8192).</i>					
Linear	8192	671K	1.56	1.61	1.56
Chain-of-Linear (Naive)	8192	419K	0.96	0.84	1.82
Chain-of-Linear (Triton)	8192	419K	0.77	0.76	0.90

A.6 Implementation

In this section, we will provide detail implementation about each module used in CoLM, including Linear, Attention, FFN, normalization and so on.

A.6.1 Linear

In Algorithm 1, we provide the naive PyTorch implementation of Chain-of-Linear layer. Generally, the proposed Chain-of-Linear layer can be viewed as a composition of multiple sub-linear layers. Under the same dimension, our implementation (i.e., Algorithm 1) can reduce more weights when compared with the standard Linear layer, but it will add the number of data loading and all-reduce communication. Therefore, we further devise a block-wise sparse kernel to accelerate the computation of Chain-of-Linear structure, which is introduced in below.

Block-wise Sparse Kernel Acceleration We developed a block-wise sparse matrix multiplication kernel for both forward and backward paths of Chain-of-Linear layer. To conserve GPU memory, we store the weights of Linear (Chain) layer in block-compressed-sparse-row (BCSR) format and perform matrix multiplication only on the step-shaped mask to accelerate computation. The kernel implementation is based on the Triton General-Matrix-Multiplication (GeMM) kernel and supports both sparse input and sparse output. For sparse input, each thread block processes a row in the BCSR matrix, skipping unmapped data from the other dense input matrix. For sparse output, thread blocks in the masked area will exit immediately. Figure 8 illustrates the details of our block-wise sparse kernel for acceleration. We also report the computation comparisons between our kernel with naive

Algorithm 2 Pseudo Code for Attention (Chain) Layer.

```
from flash_attn import flash_attn_func
from Linear_chain import Linear

class Attention(nn.Module):

    def __init__(
        self, dim: int, n_head: int, n_kv_head: int, chains: List[int]
    ):
        super().__init__()
        assert sum(chains) == n_head
        for chain in chains:
            assert chain * n_kv_head % n_head == 0
        self.head_dim = dim // n_head

        self.wq = Linear(dim, dim, chains)
        self.wk = Linear(dim, n_kv_head * self.head_dim, chains)
        self.wv = Linear(dim, n_kv_head * self.head_dim, chains)
        self.wo = Linear(dim, dim, chains)

    def forward(self,
        x: torch.Tensor,
        freqs_cis: torch.Tensor,
        chain: int = None) -> torch.Tensor:
        bsz, seqlen, dim = x.shape

        q, k, v = self.wq(x, chain), self.wk(x, chain), self.wv(x, chain)
        q = q.view(bsz, seqlen, -1, self.head_dim)
        k = k.view(bsz, seqlen, -1, self.head_dim)
        v = v.view(bsz, seqlen, -1, self.head_dim)

        q, k = apply_rotary_emb(q, k, freqs_cis)

        o = flash_attn_func(q, k, v, causal=True)
        o = o.view(bsz, seqlen, -1)
        return self.wo(o, chain=chain)
```

implementation in Table 16. From Table 16, we observe that the naive implementation is exceedingly slow in computing backpropagation gradients for the input data, while our kernel significantly reduces its computation for input data, and then improves the speed of both the forward and backward stages. Besides, our kernel can also reduce the number of all-reduce communications during the distributed training. Benefiting from this design it allows us to conduct large-scale pre-training efficiently.

A.6.2 Attention

Just as aforementioned in Section 3.2.1, we can replace all linear layer (i.e., Q, K, V, O) by using Chain-of-Linear layer. For MHA setting [7], we use $\mathcal{C} = \{c_1, \dots, c_n\}$, where $\sum_i^n c_i = h$, to determine how many heads (i.e., q_i, k_i, v_i) should be assigned to each chain. For the GQA setting [20], we require $c_i \times h_{kv}$ should be a multiple of $\text{sum}(\mathcal{C})$. For Chain-of-Attention with KV sharing mechanism, we require each c_i should be a multiple of h_{kv} , as all keys and values are computed within the first chain and then shared across the remaining chains. And then, we will replicate keys and values per chain to guarantee that their number to match the number of query. Figure 9 illustrates the differences between GQA and Attention (Chain) with the KV sharing mechanism. We also provide the PyTorch implementation of Chain-of-Attention in Algorithm 2 and its version with KV sharing in Algorithm 3.

Algorithm 3 Pseudo Code for Attention (Chain) Layer with KV sharing.

```
import torch.nn as nn

from flash_attn import flash_attn_func
from Linear_chain import Linear

class Attention(nn.Module):

    def __init__(
        self, dim: int, n_head: int, n_kv_head: int, chains: List[int]
    ):
        super().__init__()
        assert sum(chains) == n_head
        for chain in chains:
            assert chain % n_kv_head == 0
        self.head_dim = dim // n_head
        self.kv_dim = chains[0] * self.head_dim

        self.wq = Linear(dim, dim, chains)
        self.wk = nn.Linear(self.kv_dim, n_kv_head * self.head_dim, chains)
        self.wv = nn.Linear(self.kv_dim, n_kv_head * self.head_dim, chains)
        self.wo = Linear(dim, dim, chains)

    def forward(self,
                x: torch.Tensor,
                freqs_cis: torch.Tensor,
                chain: int = None) -> torch.Tensor:
        bsz, seqlen, dim = x.shape
        x0 = x[..., :self.kv_dim]

        q, k, v = self.wq(x, chain), self.wk(x0), self.wv(x0)
        q = q.view(bsz, seqlen, -1, self.head_dim)
        k = k.view(bsz, seqlen, -1, self.head_dim)
        v = v.view(bsz, seqlen, -1, self.head_dim)

        q, k = apply_rotary_emb(q, k, freqs_cis)

        n_rep = q.shape[2] // k.shape[2]
        k = k.repeat(1, 1, n_rep, 1)
        v = v.repeat(1, 1, n_rep, 1)
        o = flash_attn_func(q, k, v, causal=True)
        o = o.view(bsz, seqlen, -1)
        return self.wo(o, chain=chain)
```

A.6.3 FFN

In Figure 10, we illustrate the differences between the original FFN and our proposed FFN layer, which just needs to replace the Linear layer by using a Chain-of-Linear layer. They are equivalent when the number of chains $n = 1$. We use the same hyperparameter \mathcal{C} of the Attention module to set up the FFN layer, so that the output features of the Transformer block (i.e., Attention + FFN) also follow the criteria of CoL in Definition 2. We also offer a PyTorch implementation of FFN (GeLU) in Algorithm 4 and FFN (SwiGLU) in Algorithm 5 for reference.

Algorithm 4 Pseudo Code for FFN implementation (GeLU).

```
import torch.nn.functional as F
# Import Chain-of-Linear Layer
from Linear_chain import Linear

class FFN(nn.Module):

    def __init__(
        self, dim: int, hidden_dim: int, chains: List[int]
    ):
        """
        dim:          Dimension of input features
        hidden_dim:    Dimension of output features
        chains:        Base of each chain
        """
        super().__init__()

        self.w1 = Linear(dim, hidden_dim, chains=chains)
        self.w2 = Linear(hidden_dim, dim, chains=chains)

    def forward(self, x: torch.Tensor, chain: int = None) -> torch.Tensor:
        return self.w2(F.gelu(self.w1(x, chain)), chain)
```

Algorithm 5 Pseudo Code for FFN implementation (SwiGLU).

```
import torch.nn.functional as F
# Import Chain-of-Linear Layer
from Linear_chain import Linear

class FFN(nn.Module):

    def __init__(
        self, dim: int, hidden_dim: int, chains: List[int]
    ):
        """
        dim:          Dimension of input features
        hidden_dim:    Dimension of output features
        chains:        Base of each chain
        """
        super().__init__()

        self.w1 = Linear(dim, hidden_dim, chains)
        self.w2 = Linear(hidden_dim, dim, chains)
        self.w3 = Linear(dim, hidden_dim, chains)

    def forward(self, x: torch.Tensor, chain: int = None) -> torch.Tensor:
        return self.w2(F.silu(self.w1(x, chain)) * self.w3(x, chain), chain)
```

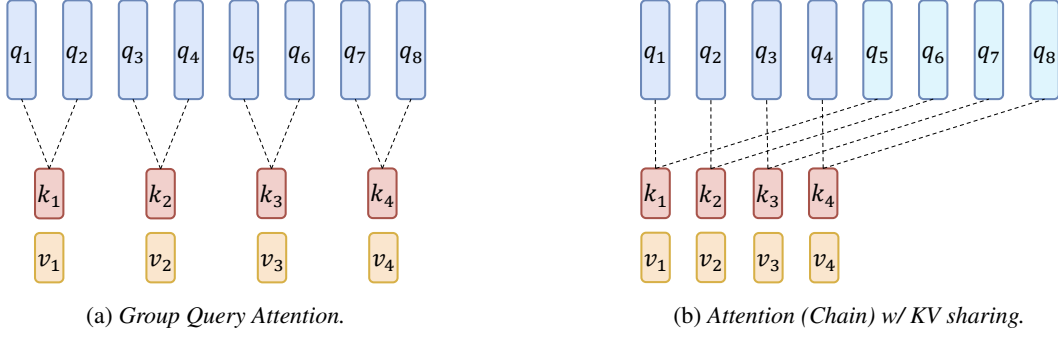


Figure 9: Differences between Group Query Attention (GQA) and Attention (Chain) with KV sharing in the attention implementation.

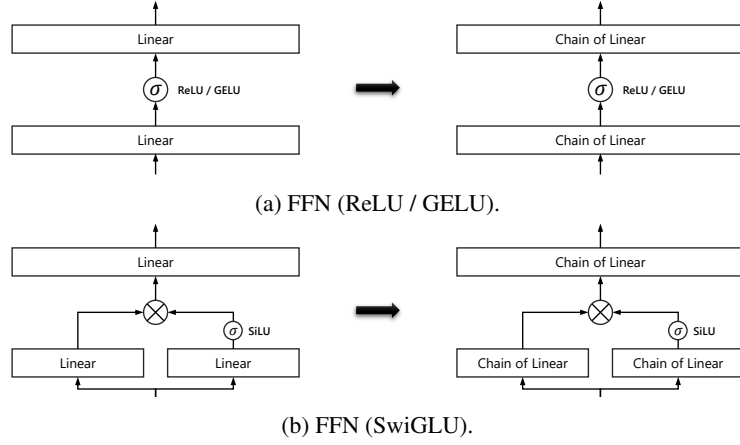


Figure 10: FFN for CoR. We replace all Linear by using Chain-of-Linear layer.

A.6.4 Embedding

Figure 11 illustrates how the Embedding layer is utilized in our CoLM architecture by activating different numbers of chains. To use different scales, we only need to activate the dimension corresponding to its all preceding chains.

Algorithm 6 Pseudo Code for Embeddings.

```
import torch.nn as nn

class Embedding(nn.Embedding):

    def __init__(
        self,
        dim: int,
        vocab_size: int,
        chains: List[int],
        padding_idx: Optional[int] = None,
    ):
        super().__init__(vocab_size, dim, padding_idx)
        self.dims = [dim * c // sum(chains) for c in np.cumsum(chains)]
        self.num_chains = len(chains)

    def forward(self, x: torch.Tensor, chain: int = None):
        if chain is None:
            return F.embedding(x, self.weight, self.padding_idx)
        else:
            assert chain < self.num_chains
            dim = self.dims[chain]
            return F.embedding(x, self.weight[:, :dim], self.padding_idx)
```

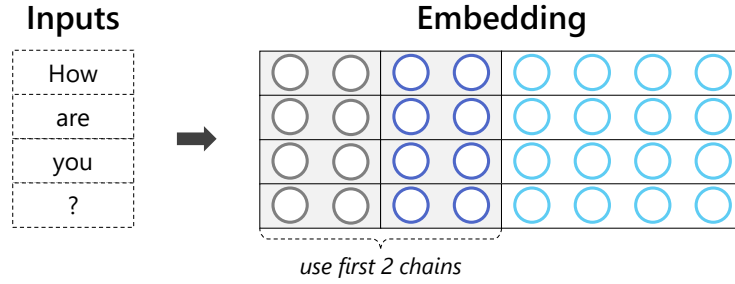


Figure 11: Example of Embedding in CoLM architecture. In this example, we design the number of chains as 3, while the dimensions for each chain are set as $\{2, 2, 4\}$. If we expect to use the first 2 chains, we only need to activate the first 4 neurons for each word in the embedding layer.

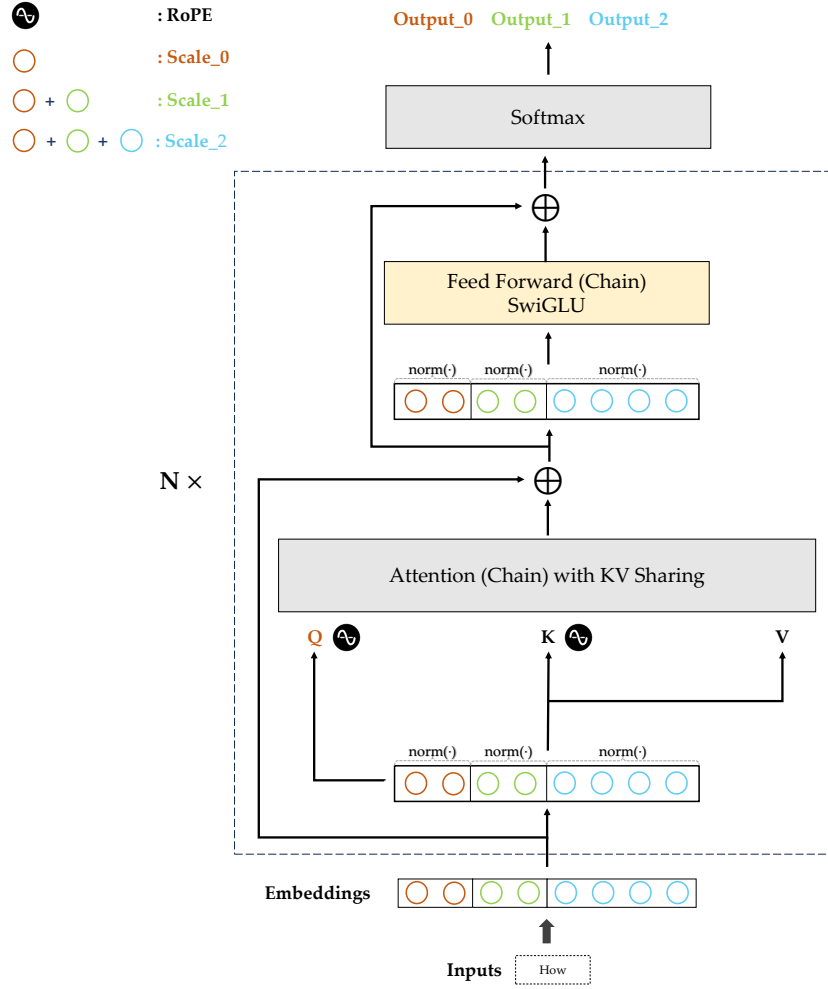


Figure 12: The architecture of our CoLM. The model use Chain-of-Representation to build each layer, including embedding, self-attention, feed-forward network, normalization. Each scale contributes progressively refined outputs (e.g., Output0, Output1, Output2).

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: The abstract and introduction clearly state the paper's contributions, scope, and assumptions, which are consistent with the presented results.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: Yes. See Section Limitation for details.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[NA\]](#)

Justification: This paper does not include theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: The model architecture and training procedure are fully specified in Section 4, enabling reproduction.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: Yes. The code and instructions for reproducing the main experimental results are publicly available on GitHub (link provided in the abstract).

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Yes. Training details, hyperparameters, and data splits are described in Section 4 and Appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: No. Due to the high computational cost, experiments were run with a single seed.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)

- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [\[Yes\]](#)

Justification: Yes. Compute resources and training time are reported in Section 4.1.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [\[Yes\]](#)

Justification: Yes. We have read the NeurIPS Code of Ethics and confirm that this research complies with it.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [\[NA\]](#)

Justification: This work is primarily methodological and foundational, without a direct or immediate path to real-world deployment or identifiable societal impacts.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper does not release models or datasets that pose a high risk for misuse.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All existing code, data, and models used in this work are properly cited in the paper, and their licenses and terms of use are respected.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. **New assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [\[Yes\]](#)

Justification: The released code is documented in the GitHub repository, including usage instructions and license information.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and research with human subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [\[NA\]](#)

Justification: This work does not involve crowdsourcing or research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [\[NA\]](#)

Justification: This work does not involve crowdsourcing or research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: LLMs were only used for writing and editing purposes and did not affect the core methodology or results.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>) for what should or should not be described.