COMPRESSING LARGE MOE MODELS VIA EFFICIENT PRUNING AND DATA-AWARE CALIBRATION

Anonymous authors

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

033

036

038

040

041

042

043

044

045

046 047

048

051

052

Paper under double-blind review

ABSTRACT

Ultra-large Mixture-of-Experts (MoE) language models, e.g., DeepSeek-R1, are rapidly emerging as a dominant architecture due to their superior scalability and performance. However, the massive number of expert parameters introduces substantial redundancy, posing serious challenges for efficient deployment. Existing pruning methods face two fundamental challenges when applied to such MoE architectures. First, while methods based on reconstruction loss offer a more comprehensive selection by considering each expert combination, the vast search space renders exhaustive evaluation infeasible. Second, most approaches rely on a fixed calibration dataset to guide pruning, which often fails to preserve the model's full capabilities. To address these challenges, we introduce two key innovations in our pruning framework. First, we propose a Coarseto-Fine Expert Selection strategy that reduces the computational complexity of reconstruction-loss-based selection from an exponential $(\mathcal{O}(\binom{2n}{n}))$ to a polynomial scale $(\mathcal{O}(n^{1.5}))$ with respect to the number of experts. This significantly accelerates the pruning process without sacrificing selection quality. Second, we develop a Dynamic Calibration Dataset Mixing strategy that enables the model to adaptively adjust its calibration set during pruning. Extensive experiments on a range of benchmarks show that our method can prune 50% of the experts in a large-scale MoE model (e.g., DeepSeek-R1) while retaining 98.9% of its original performance across diverse tasks, outperforming existing pruning baselines. Our approach also demonstrates practical speedups and reduced memory footprint, facilitating efficient real-world deployment. The anonymous implementation is available at https://anonymous.4open.science/r/ DCDM-4C65-622a2bad88498795b8d7a92d85aca1315f9520ee.

1 Introduction

Mixture-of-Experts (MoE) models (Zhao et al., 2023; Team, 2024; Jiang et al., 2024; Dai et al., 2024a; Fedus et al., 2022) have demonstrated remarkable performance across a wide range of natural language processing tasks due to their ability to scale model capacity via a large number of experts while keeping per-token computation relatively constant by activating only a small subset at each step. However, a major bottleneck in deploying MoE models is the memory overhead introduced by inactive experts (Gao et al., 2022), which must still be stored during inference despite not contributing to computation. This issue is particularly severe for ultra-large models with a massive pool of experts (DeepSeek-AI et al., 2025; Team, 2025), of which only a small fraction are ever activated. If inactive experts could be effectively removed or reduced, the memory and storage costs would drop significantly, making MoE models much more practical for real-world deployment.

A straightforward and widely used approach is to prune experts based on router-derived metrics, such as gate values or activation frequency (Cao et al., 2024b; Dong et al., 2025). This strategy is computationally inexpensive and requires only a single forward pass, making it attractive for large-scale models. However, such metrics treat each expert independently and fail to account for the mutual influence among experts, often leading to suboptimal pruning decisions. To more accurately measure an expert's contribution, some work introduces reconstruction-loss-based selection (Lu et al., 2024), which evaluates the change in model outputs when an expert is removed. Although this provides a more principled criterion, it typically requires evaluating a large number of

expert combinations, which becomes computationally prohibitive as the number of experts grows. To further reduce redundancy while preserving knowledge, clustering-based methods (Chen et al., 2025) attempt to merge similar experts, yet these approaches still incur high computational cost and degrade performance in ultra-large models. Despite these advances, designing a pruning strategy that balances accurate expert contribution estimation with computational efficiency remains a key challenge that we aim to address in this work.

To address these challenges, we introduce a pruning framework for MoE models grounded in reconstruction-loss-based selection. First, to alleviate the massive computation induced by huge number of experts, we adopt a *layer-wise greedy search* that incrementally selects critical experts in each layer by minimizing the discrepancy between the outputs of the original and pruned layers. This design reduces the complexity from exponential, $\mathcal{O}(\binom{2n}{n})$, to polynomial, $\mathcal{O}(n^2)$, and can be further accelerated by a *coarse-to-fine expert selection mechanism* that lowers the cost to $\mathcal{O}(n^{1.5})$, enabling pruning even for ultra-large MoE models. We also provide a theoretical guarantee that the global error can be bounded by the accumulated layer-wise error, ensuring alignment with the overall objective. Second, to overcome the reliance on domain-specific calibration datasets, we introduce a *dynamic calibration dataset mixing* strategy that adaptively adjusts the mixture of samples from different domains based on the discrepancy between the original and pruned models, thereby enhancing the generalization of the pruned model across diverse domains. Together, these two innovations enable scalable and generalizable pruning for MoE models, making them far more practical for real-world deployment.

Extensive experiments validate the effectiveness of our framework. We evaluate the pruned models on math (*i.e.*, AIME) and code (*i.e.*, LiveCodeBench) benchmark, our pruned model achieves up to 98.9% of the original model's performance while reducing memory usage by 50%. Notably, it remains competitive even on challenging multi-domain scenarios, where existing pruning methods often suffer significant degradation. These results highlight that our approach not only scales to ultralarge MoE models but also preserves generalization across diverse domains, making it significantly more deployment-friendly in real-world settings.

2 Related work

In this section, we first trace the evolution of MoE architectures that enable large-scale language models to scale efficiently via sparse expert activation. Second, we examine the major bottleneck in MoE deployment, the memory overhead caused by inactive experts, and review recent efforts that aim to prune expert parameters to reduce deployment cost.

2.1 EFFICIENT MOE ARCHITECTURES

The introduction of sparsely activated Mixture-of-Experts (MoE) architectures has been a key driver in scaling large language models. Early works such as GShard (Lepikhin et al., 2021) and the Switch Transformer (Fedus et al., 2022) demonstrated that activating only a small subset of experts per token can dramatically increase model capacity without proportionally increasing computation. Building on these foundations, the BASE Layers (Lewis et al., 2021) formulation modeled token–expert assignments as a constrained optimization problem, enabling balanced expert usage at scale. Meanwhile, Task-MoE (Kudugunta et al., 2021) extended the routing paradigm by assigning experts at the task level, reducing routing variance and improving efficiency in multilingual settings. More recent systems, including DeepSeek-MoE (Dai et al., 2024b) and Mixtral (Jiang et al., 2024), explored hierarchical routing and large-scale deployment settings, further improving throughput and stability. These architectural innovations collectively highlight the efficiency gains achieved through sparse expert activation, laying the groundwork for subsequent research on reducing redundancy among experts and optimizing their deployment. These designs typically require training from scratch. In contrast, a more practical challenge is how to reduce the number of experts without reinitializing or retraining the full model.

2.2 Compression of MoE models

Despite the computational advantages of sparse activation, MoE models still suffer from high memory and storage overhead due to the large number of inactive experts. To address this issue, a growing

line of research has investigated expert pruning as a means of compressing MoE models while maintaining performance. Existing methods further fall into three categories. The first category leverages routing-based metrics(Lu et al., 2024; Chen et al., 2025) to identify important experts. These approaches are computationally efficient but typically evaluate experts independently, without assessing their global impact on model behavior. The second category employs reconstruction-loss-based strategies (Cao et al., 2024a), which measure the discrepancy between outputs of the original and pruned models, providing a more principled criterion by directly minimizing changes in functional behavior. The third category merges experts through clustering algorithms (Chen et al., 2025), aiming to preserve knowledge from all experts, though such methods often underperform when scaling to very large models. In this work, we follow the reconstruction-loss-based direction, formalizing expert pruning as an output discrepancy minimization problem.

3 PRELIMINARIES

Notations and objective. In general, we consider a Mixture-of-Experts (MoE) model with L layers. Let X_l and X_{l+1} denote the input and output of the l-th layer, respectively. Each Layer $_l$ contains a router function g_l that selects a subset of experts (e.g., top-8 in DeepSeek-R1) from a total of c experts. The full expert set in the l-th layer is denoted by $\mathcal{E}_{l,c}$, and any subset with \hat{c} experts is written as $\mathcal{E}_{l,\hat{c}}$. The output of the l-th layer is given by $X_{l+1} = \mathrm{Layer}_l(X_l; \mathcal{E}_{l,c})$. For notational brevity, we may simplify $\mathrm{Layer}_l(\cdot; \mathcal{E}_{l,c})$ as $\mathrm{Layer}_{l,c}$ when the input and expert set are clear from context. Then the overall model output is defined as:

$$X_{L+2} = \left(\operatorname{Layer}_{L+1} \circ \operatorname{Layer}_{L,c} \circ \cdots \circ \operatorname{Layer}_{1,c}\right)(X_1), \tag{1}$$

where $\operatorname{Layer}_{L+1}$ denotes the head layer. Let \hat{X}_l and \hat{X}_{l+1} be the input and output of the l-th pruned layer, then the output of the pruned model at the l-th layer is given by $\hat{X}_{l+1} = \operatorname{Layer}_l(\hat{X}_l; \mathcal{E}_{l,\hat{c}}) := \operatorname{Layer}_{l,\hat{c}}(\hat{X}_l)$, with $\hat{X}_1 = X_1$ as the input to the first MoE layer. The goal of pruning is to obtain, for each layer, a reduced expert subset $\mathcal{E}_{l,\hat{c}} \subset \mathcal{E}_{l,c}$ such that the output of the pruned model:

$$\hat{X}_{L+2} = \left(\text{Layer}_{L+1} \circ \text{Layer}_{L,\hat{c}} \circ \dots \circ \text{Layer}_{1,\hat{c}} \right) (X_1), \tag{2}$$

approximates X_{L+2} as closely as possible. Let $d(\cdot, \cdot)$ be a distance metric measuring the discrepancy between X_{L+2} and \hat{X}_{L+2} . The pruning objective can then be formulated as:

$$\min_{\mathcal{E}_{l,\hat{c}}} d(X_{L+2}, \hat{X}_{L+2}), \quad \text{s.t.} \quad \mathcal{E}_{l,\hat{c}} \subset \mathcal{E}_{l,c}. \tag{3}$$

With this notation, we formalize expert pruning as selecting $\mathcal{E}_{l,\hat{c}}$ to minimize the output discrepancy. We next introduce our method for solving this problem efficiently.

4 METHODS

This section introduces our pruning framework, which primarily focuses on efficient expert selection and also incorporates a dynamic calibration strategy to improve cross-domain generalization. First, we present an expert selection method grounded in output discrepancy. We prove that the global pruning discrepancy can be bounded by the cumulative layer-wise discrepancy, providing a theoretical justification for adopting layer-wise pruning. Building on this analysis, we design a coarse-to-fine expert selection strategy (left side of Figure 1). Additionally, we introduce a dynamic calibration dataset mixing strategy (right side of Figure 1), which adaptively blends domain-specific and general data to further enhance the pruned model's generalization.

4.1 EXPERT SELECTION BASED ON DISCREPANCY

Layer-wise greedy search. The objective of pruning is to ensure that the output of the pruned model closely approximates that of the original model. However, directly optimizing expert selection across all layers to minimize the final output discrepancy is computationally impractical, as it would require full forward passes through the model at every search step. A simpler alternative is to instead minimize the discrepancy between the original and pruned outputs at each layer individually. This raises a natural question:

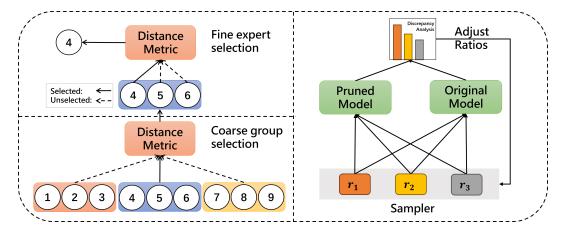


Figure 1: Overview of the proposed pruning framework. **Left**: Coarse-to-fine expert selection first scores groups of experts by output discrepancy, then refines the choice within the selected group. **Right**: Dynamic calibration dataset mixing adjusts domain sampling according to pruning-induced discrepancy, enhancing generalization across domains.

Can the overall model output discrepancy be effectively controlled by minimizing layer-wise discrepancy?

To address this, we propose the following theorem, which shows that the global output difference can indeed be bounded by the accumulation of local layer-wise discrepancy.

Theorem 4.1 (Layer-wise Pruning Bound). Let X_{L+2} denote the output of the original model, \hat{X}_{L+2} denote the output of the pruned model and $d(\cdot, \cdot)$ be a distance metric. Suppose each layer is pruned from c experts to \hat{c} experts. Then the distance $d(X_{L+2}, \hat{X}_{L+2})$ between the original and pruned outputs is bounded as follows:

$$d(X_{L+2}, \hat{X}_{L+2}) \le \sum_{i=2}^{L+1} \operatorname{Lip}_{i \to L+1} d(\hat{X}_i, \operatorname{Layer}_{i-1, c}(\hat{X}_{i-1}))$$
(4)

where $Lip_{i\rightarrow L+1}$ denotes the Lipschitz constant from the i-th layer to the (L+1)-th layer.

The complete proof is provided in Appendix B.

Theorem 4.1 shows that the discrepancy between the outputs of the original and the pruned model can be bounded by the accumulation of layer-wise differences. This result provides a theoretical justification for replacing the computationally expensive full-model comparison $d(X_{L+2}, \hat{X}_{L+2})$ with layer-wise comparisons $d(\hat{X}_i, \text{Layer}_{i-1,c}(\hat{X}_{i-1}))$. Then, the costly process of repeatedly performing full forward passes and evaluating combinations of experts across multiple layers can be avoided, and the search space is reduced to experts within a single layer.

However, exhaustively evaluating all expert combinations within one layer still remains computationally expensive, exhibiting exponential time complexity (Lu et al., 2024) for layers with many experts (e.g., 256 experts in DeepSeek-R1). Inspired by prior work (Cao et al., 2025), we adopt a greedy strategy that selects the most critical expert from each layer—the one that minimizes the output difference between the pruned and original models at a search time. This greedy approach reduces the time complexity from exponential to polynomial. Although polynomial complexity is more manageable, the computational cost can still be high during the search process with a large number of experts. To address this, we next introduce a coarse-to-fine expert selection strategy to further reduce the time cost.

Coarse-to-fine expert selection. To accelerate the otherwise expensive greedy search over all candidate experts, we adopt a coarse-to-fine strategy (illustrated on the left of Figure 1). The key idea is to first evaluate experts at a group level to quickly narrow down promising candidates, and then

241242243244245246247248249

250

253

254

255

256257

258

259

260

261

262

264

265

266

267 268

```
216
                Algorithm 1 Coarse-to-Fine Expert Selection
217
218
                 1: Calibration dataset: X_1, Selected experts: \hat{c}.
219
                 2: Architecture: 1)Total layers L. 2) Experts per layer N_e.
220
                 3: Experts initialization:
221
                            • Selected experts at layer l: \mathcal{E}_{l,\text{selected}} = \emptyset (\forall l \in \{1, \dots, L\}).
                            • Candidate experts at layer l: \mathcal{E}_{l, \text{candidate}} = \{1, 2, \dots, N_e\} (\forall l \in \{1, \dots, L\}).
222
                            • Group size: S = \text{round}(\sqrt{N_e + \frac{1-\hat{c}}{2}}).
224
                Ensure: Pruned expert set \{\mathcal{E}_{l,\text{selected}}\}_{l\in\{1,\ldots,L\}}.
225
                 4: Compute each layer output of X_{l+1} for each Layer, via calibration dataset X_1
226
                  5: for each MoE layer l to L do
227
                 6:
                             \mathcal{E}_{\text{selected}} \leftarrow \mathcal{E}_{l, \text{selected}}, \mathcal{E}_{\text{candidate}} \leftarrow \mathcal{E}_{l, \text{candidate}}
228
                             K \leftarrow \operatorname{ceil}(|\mathcal{E}_{\operatorname{candidate}}|/S)
                 7:
                 8:
                             for t=1 to \hat{c} do
229
                 9:
                                   Partition \mathcal{E}_{\text{candidate}} into K groups \{\mathcal{G}_1, ..., \mathcal{G}_K\}
230
                10:
                                   for each group G_k do
231
                11:
                                         \mathcal{E}_{\text{temp}} \leftarrow \mathcal{E}_{\text{selected}} \cup \mathcal{G}_k
232
                                         Compute output distance d_k \leftarrow d(X_{l+1}, \text{Layer}_l(X_l; \mathcal{E}_{\text{temp}}))
                12:
                13:
                                    Find optimal group \mathcal{G}^* \leftarrow \arg\min_{\mathcal{G}_k} d_k
                15:
                                    Extract best expert e^* \leftarrow \arg\min_{e \in \mathcal{G}^*} d(\operatorname{Layer}_l(X_l; \mathcal{E}_{\text{selected}} \cup \{e\}), X_{l+1})
235
                                    Update: \mathcal{E}_{\text{selected}} \leftarrow \mathcal{E}_{\text{selected}} \cup \{e^*\}
                16:
                17:
                                   Update: \mathcal{E}_{\text{candidate}} \leftarrow \mathcal{E}_{\text{candidate}} \setminus \{e^*\}
237
                18:
                             end for
238
                19:
                             \mathcal{E}_{l,\text{selected}} \leftarrow \mathcal{E}_{\text{selected}}
239
               20: end for
```

Table 1: Analysis of the time complexity.

Methods	Time Complexity	Time Complexity in Terms of n
O-Prune (Lu et al., 2024)	$\mathcal{O}\left(FL\binom{N_e}{\hat{c}}\right)$	$\mathcal{O}({2n \choose n}) \ \mathcal{O}(n^3)$
HC-SMoE (Chen et al., 2025)	$\mathcal{O}(L(N_e - \hat{c})N_e^2 + FL)$	$\mathcal{O}(n^3)$
Layer-wise greedy search	$\mathcal{O}(\hat{c}FL(N_e + \frac{1-\hat{c}}{2}))$	$\mathcal{O}(n^2)$
Coarse-to-fine expert selection	$\mathcal{O}(\hat{c}FL\sqrt{4N_e+2-2\hat{c}})$	$\mathcal{O}(n^{1.5})$

refine the choice within the selected group using a finer-grained metric. Concretely, the candidate expert set $\mathcal{E}_{l,c}$ at layer l is first partitioned into K disjoint groups $\mathcal{G}_1,\ldots,\mathcal{G}_K$, where the first K-1 groups contain S experts each and the last group holds the remaining candidates. The selection proceeds in two stages:

- Coarse group selection: For each group \mathcal{G}_i , we compute the output discrepancy between the original layer and a pruned layer that uses the currently selected experts $\mathcal{E}_{l,\text{selected}}$ together with all experts in \mathcal{G}_i . A distance metric (the ℓ_2 -norm in our experiments) measures this discrepancy. The group \mathcal{G}^* with the smallest discrepancy is chosen as the most promising region.
- Fine expert selection: Within \mathcal{G}^* , we further evaluate each expert individually using the same distance metric and select the expert e^* that yields the closest match to the original layer output when added to $\mathcal{E}_{l,\text{selected}}$.

The complete expert pruning algorithm is formally presented in Algorithm 1.

Analysis of the time complexity. We analyze the computational cost of our coarse-to-fine expert selection strategy in this part. Generally, we denote each MoE layer route to N_e candidate experts and require F time for a single forward pass. Then a naive greedy selection evaluates every remaining expert at each iteration, resulting in a large search space. Instead, our hierarchical strategy significantly shrinks this space. At the t-th iteration, the candidate pool is first reduced from $|N_e-t+1|$ experts to $\lceil \frac{N_e-t+1}{S} \rceil + S$. The resulting overall complexity is:

$$\mathcal{O}((\sum_{t=1}^{\hat{c}} \lceil \frac{N_e - t + 1}{S} \rceil + \hat{c}S)FL) = \mathcal{O}(\hat{c}FL(S + \frac{2N_e + 1 - \hat{c}}{2S})), \tag{5}$$

Algorithm 2 Dynamic Calibration Dataset Mixing Strategy

```
271
             Require: Initial weights \mathbf{w}^0, the required in Algorithm 1
272
              1: for pruning iteration t = 1 to T do
273
                       Sample \mathcal{X}_1 from \{D_1, \cdots, D_{N_D}\} with proportion \boldsymbol{w}^{t-1}.
274
                       Prune model through Algorithm 1
275
              4:
                       Evaluate the discrepancy d^t across N_D domains
                       Update weights: \mathbf{w}^t \leftarrow \frac{\exp \boldsymbol{d}^t}{\|\exp \boldsymbol{d}^t\|_1}
276
              5:
                       if \mathbf{w}^t == \mathbf{w}^{t-1} then break
277
278
             7:
                       end if
             8: end for
279
```

where the term $S+\frac{2N_e+1-\hat{c}}{2S}$ captures the trade-off between coarse grouping and fine expert selection. Then the optimal time complexity $\mathcal{O}(\hat{c}FL\sqrt{4N_e+2-2\hat{c}})$ is achieved when the group size $S=\sqrt{N_e+\frac{1-\hat{c}}{2}}$, which is a significant improvement over the complexity of a naive greedy selection. We compare the time complexity of different methods in Table 1. By setting $N_e=2n$ and $\hat{c}=n$, while treating other parameters as constants, the time complexity can be asymptotically expressed in terms of n.

To quantitatively assess the computational efficiency of our method relative to existing approaches, we examine the DeepSeek-R1 model as an example. It contains L=58 layers with $N_e=256$ routed experts per layer. When selecting $\hat{c}=128$ experts, the method of O-Prune (Lu et al., 2024) would require $58\binom{256}{128}$ evaluations, HC-SMoE (Chen et al., 2025) about 4.87×10^8 , and a layer-wise greedy search about 1.43×10^6 . In contrast, our method needs only 2.06×10^5 evaluations, showing a substantial reduction in time cost and markedly better efficiency than competing approaches.

Overall, our coarse-to-fine expert selection scheme reduces computational cost and provides a practical approach for pruning large language models with numerous experts within a reasonable time.

4.2 Dynamic Calibration dataset mixing strategy

While existing work shows that domain-specific performance can often be well preserved, expert pruning still suffers from notable cross-domain degradation (Dong et al., 2025). To alleviate the domain shift introduced by the calibration data, we propose a Dynamic Calibration Dataset Mixing (DCDM) strategy, motivated by recent work (Xie et al., 2023; Xia et al., 2024).

Specifically, given N_d distinct domains $\{\mathcal{D}_i\}_{i=1}^{N_D}$, we initialize the domain mixing weights $\mathbf{w}^0 = (w_1^0,\dots,w_{N_D}^0)$ proportionally to their dataset sizes $w_i^0 = \frac{|\mathcal{D}_i|}{\sum_{j=1}^{N_D}|\mathcal{D}_j|}$. Let $\mathbf{d}^t = (d_1^t,\dots,d_{N_D}^t)$ denote the discrepancy between the original model and the pruned model outputs across N_D domains at the t-th pruning iteration.

Accordingly, our objective is to minimize the performance gap between the pruned and original models, ensuring the outputs of the pruned model remains as close as possible to that of the original model. Specifically, if the discrepancy of domain \mathcal{D}_i is larger, then increase the sample ratio of the \mathcal{D}_i domain, otherwise, it is decreased. Formally, at the t-th pruning iteration the ratio is obtained by $\mathbf{w}^t = \frac{\exp(\mathbf{d}^t)}{\|\exp(\mathbf{d}^t)\|_1}$. And the dynamic calibration dataset mixing strategy is presented in Algorithm 2.

5 EXPERIMENTS

5.1 EXPERIMENT SETUP

Base models. We conduct expert pruning on two popular large MoE models: DeepSeek-R1 and Qwen3-30B-A3B-Thinking (Team, 2025).

• **DeepSeek-R1.** It is a large-scale MoE model with 671B total parameters. It contains 61 layers, including 3 dense transformer layers and 58 MoE layers. Each MoE layer has 256 routed experts and 1 shared expert. During inference, the router selects 8 routed experts per layer, which are combined with the shared expert for computation.

Table 2: Performance comparison with expert pruning on Qwen3-30B-A3B-Thinking.

Dataset	Method	MMLU	Math500	AIME25	LCB	Average
-	Original	78.59	96.12	85.00	66.00	81.43
	Weights	46.99	14.57	0.42	0.00	15.50
C4	HC-SMoE	39.42	72.98	22.08	3.18	34.41
	Ours	59.91	3.73	0.42	0.00	16.02
	Weights	33.36	13.65	2.50	4.90	13.60
OpenR1-Math	HC-SMoE	47.06	73.20	20.00	32.00	43.07
	Ours	44.10	94.60	73.75	0.60	53.26
	Weights	39.31	20.30	4.16	0.00	15.94
rStar-Coder	HC-SMoE	47.05	72.08	24.58	2.40	36.53
	Ours	44.54	93.08	71.67	63.00	68.07
Mixed Datasets	DCDM	52.56	95.60	80.00	52.40	70.14

• Qwen3-30B-A3B-Thinking. It is a 30.5B parameter MoE model with 48 MoE layers. Each layer contains 128 routed experts, with 8 experts selected per token during inference.

In our experiments, we apply a 50% sparsity ratio to the routed experts, reducing the number of experts per MoE layer from 256 to 128 for DeepSeek-R1, and from 128 to 64 for Qwen3-30B-A3B-Thinking.

Baseline methods. To evaluate the effectiveness of our method, we compare against two representative baselines chosen along two key dimensions: the pruning criterion (static statistics *vs.* output-aware perturbation) and the pruning strategy (expert pruning *vs.* expert merging). To enable a fair comparison, we carefully select baseline methods by varying one factor at a time, either the pruning criterion or the pruning mechanism, while keeping other variables such as the underlying model architecture and sparsity ratio consistent.

- Weights (Routing-based: Following CD-MoE (Cao et al., 2024b), experts are ranked by the ratio of average routing weight to activation frequency. This favors experts that are both frequently used and strongly weighted, mitigating bias toward rarely-used or weakly-activated ones. The top- \hat{c} experts are retained.
- HC-SMoE (Clustering-based): HC-SMoE (Chen et al., 2025) merges semantically similar experts based on their impact on model outputs, rather than pruning them. This clustering-based strategy preserves output quality while reducing the expert count.

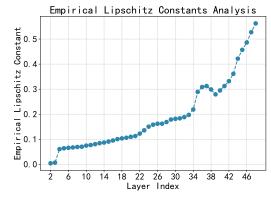
Datasets. Calibration data during pruning are sourced from three distinct task domains: C4 (Raffel et al., 2019) (knowledge), OpenR1-Math (LI et al., 2024; open r1, 2025) (mathematical reasoning), and rStar-Coder (Liu et al., 2025) (code synthesis). As for the calibration size, we use 32 calibration examples of approximately 4K tokens each under single-domain settings. In our dynamic calibration dataset mixing strategy, 32 examples are selected in total, with the domain mixture automatically adjusted based on the strategy's feedback. The pruned model's performance is evaluated on four established benchmarks: MMLU (Hendrycks et al., 2020) (knowledge), Math500 (Lightman et al., 2023) and AIME25 (math ai, 2025) (mathematical reasoning), and LiveCodeBench (Jain et al., 2024) (abbreviated "LCB"; code synthesis).

5.2 Main Results

Pruning results on Qwen3-30B-A3B-Thinking. We evaluate our approach on the Qwen3-30B-A3B-Thinking model to analyze the effect of different calibration strategies, as summarized in Table 2. We first examine pruning with a calibration dataset drawn from a single domain. This setting often achieves strong in-domain performance but suffers from severe cross-domain degradation. For example, pruning with C4 data yields competitive results on MMLU (*i.e.*, 59.67) but fails on mathematical and coding tasks (*e.g.*, 0.42 on AIME25). Similarly, pruning with domain-specific data such as OpenR1-Math excels in mathematics (*e.g.*, 94.60 on Math500) but generalizes poorly to other tasks (*e.g.*, 0.60 on LCB). These observations highlight a key limitation: single-domain pruning tends to overfit the calibration domain, causing catastrophic performance drops in other areas. To overcome this issue, we introduce the proposed Dynamic Calibration

Table 3: Comparison of Data Mixing Strategies.

Method	MMLU	Math500	AIME25	LCB	Average
Fixed	50.58	93.30	70.00	19.50	58.35
DCDM	52.56	95.60	80.00	52.40	70.14



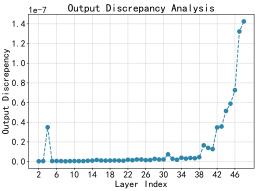


Figure 2: Empirical Lipschitz constants and output discrepancy versus layer index. (a) Lipschitz constants from the i-th layer to the 48th layer. (b) L2-norm between outputs of pruned and original models at each layer.

Dataset Mixing (DCDM) strategy, which adaptively reweights calibration data across domains according to pruning-induced output discrepancy. As shown in the final row of Table 2, DCDM achieves the highest average score (70.47), retaining 86.14% of the original model's performance (81.43) while outperforming all baselines. Notably, it avoids the severe performance collapses observed in single-domain pruning and delivers substantially better robustness and cross-domain generalization.

5.3 EXPERIMENTAL ANALYSIS

Empirical Analysis of Layer-wise Pruning Bound. To empirically assess the tightness of the upper bound derived in Theorem 4.1, we estimate the empirical Lipschitz constants $Lip_{i\rightarrow L+1}$ for the Qwen3-30B-A3B-Thining model, which comprises 48 MoE layers. As shown in the left of Figure 2, most of estimated constants are below 0.5. This indicates that the coefficients $Lip_{i\rightarrow L+1}$ multiplying the layer-wise discrepancy in our theoretical bound are small across layers. The right side shows shows correspondingly small output discrepancies (most $\leq 1.4 \times 10^{-7}$) at each layer. These results provide empirical evidence for the tightness of our theoretical bound and support the effectiveness of the layer-wise pruning strategy.

Time Cost vs. Number of Experts. As shown on the left side of Figure 3, we measure the time cost per layer as the number of experts increases (128, 256, 512, and 1024). The plotted points represent actual measured costs, while the dashed line shows a fitted curve. The trend clearly demonstrates that the time cost of our method grows at a slower rate compared to other methods as the number of experts increases. This scaling behavior aligns with our time complexity analysis presented in Table 1, confirming that our approach remains efficient and practical for modern MoE models with large numbers of experts.

Time Cost vs. Group Size. The group size S is a key hyper-parameter that balances the costs of the coarse and fine selection stages. We analyze its impact on the time cost for a fixed model size ($N_e=128$ experts pruned to $\hat{c}=64$), with results shown on the right side of Figure 3. Theoretically, the time complexity first decreases and then increases with S. Our empirical results strongly support this theoretical analysis. The measured time cost exhibits a characteristic U-shaped curve. The cost initially decreases sharply as S increases from small values. It reaches a minimum around S=10, which aligns closely with the theoretically predicted optimal value of $\operatorname{round}(\sqrt{128+(1-64)/2})=10$. Beyond this point, the cost increases gradually with further

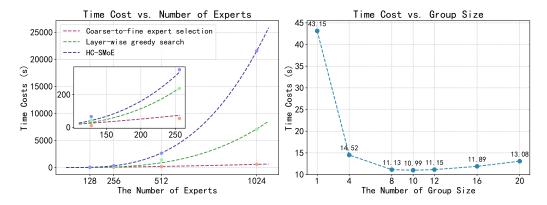


Figure 3: Time costs vs. number of experts and group size. Left: time cost versus the number of experts, with the inset providing a detailed perspective from 64 to 256 experts. Right: time cost versus group size.

increases in S. The close agreement between theory and experiment demonstrates that our time complexity analysis is correct and identifies an optimal group size for minimizing pruning time.

Data Mixing Strategies. To evaluate our dynamic calibration dataset mixing (DCDM) strategy, we compare it with a static policy that combines rStar-Coder, OpenR1-Math, and C4 data in a fixed 1:1:1 ratio for the calibration set. As shown in Table 3, the fixed strategy performs reasonably well on math-related tasks (*e.g.*, 93.30 on Math500) but struggles on code-related tasks (19.50 on LCB). This indicates that a fixed ratio may not reflect the actual data needs of each domain, for instance, code tasks may require a higher proportion of domain-specific calibration samples. However, such requirements are typically unknown prior to pruning. Our DCDM strategy adjusts the calibration set composition based on pruning feedback, resulting in a modified ratio of 2:1:1. This adaptation increases the proportion of code domain data, which was underperforming with the 1:1:1 ratio. Consequently, our method improves performance across all domains, particularly on code tasks (19.50 vs. 58.10 on LCB). These results highlight the importance of calibration data composition in expert pruning and demonstrate that our dynamic approach can effectively balance data needs across domains to preserve overall model performance.

Impact on Model Size. To further evaluate the effectiveness and scalability of our approach on larger MoE models, we conduct experiments on the representative ultra-large model DeepSeek-R1, tested on two challenging benchmarks: AIME25 for math and LCB for code. The results, summarized in Table 4, show that the pruned models retain around 98.9% of the original model's performance while significantly reducing the number of ac-

Table 4: Performance on DeepSeek-R1.

Method	AIME25	LCB	Average
Original	65.00	59.14	73.81
Ours	62.50	60.22	72.64

tive experts. Interestingly, we observe that larger models tend to maintain a higher fraction of their original performance after pruning, suggesting that scale improves robustness to compression. This observation is consistent with prior findings (Liu et al., 2024) that larger language models are generally more resilient to parameter reduction. These results highlight the scalability of our method and suggest promising potential for applying it to even larger future MoE models.

6 Conclusion

In this work, we propose an efficient expert pruning framework for large-scale Mixture-of-Experts (MoE) models, focusing on both pruning efficiency and performance preservation. Specifically, we introduce a coarse-to-fine polynomial selection strategy that reduces the search complexity from exponential scale to polynomial scale and a dynamic calibration data mixing strategy that adaptively adjusts calibration samples to improve the generalization among different domains. Experiments across diverse domains and model scales show that our method surpasses existing baselines, achieving high compression rates while retaining up to 98.9% of the original model's performance, demonstrating its practicality and scalability for real-world MoE deployment.

ETHICS STATEMENT

This work focuses on reducing memory cost of large-scale Mixture-of-Experts (MoE) models through expert pruning. Our research does not involve human subjects, personal data, or sensitive information, and all datasets used are publicly available and widely adopted in the community. We employed large language models only for grammatical checking and minor language polishing of the manuscript; they were not involved in designing experiments or generating results. We believe there is no foreseeable ethical risk beyond the general considerations of deploying more efficient large language models.

REPRODUCIBILITY STATEMENT

We provide an anonymous implementation of our method to facilitate reproduction of all reported results. A detailed explanation of the proposed method and the proofs of theoretical justification is included in the Appendix B. The datasets and experimental settings are thoroughly described in Section 5.1, including datasets, baselines, and evaluation benchmarks, to ensure that readers can replicate our experiments without ambiguity. All hyperparameters and training configurations can be found in Section 5.1. At last, we used the public available pretrained models (*e.g.*, DeepSeek-R1 and Qwen3-30B-A3B-Thinking) for easy reproducibility.

REFERENCES

- Mingyu Cao, Gen Li, Jie Ji, Jiaqi Zhang, Xiaolong Ma, Shiwei Liu, and Lu Yin. Condense, don't just prune: Enhancing efficiency and performance in moe layer pruning. *CoRR*, abs/2412.00069, 2024a.
- Mingyu Cao, Gen Li, Jie Ji, Jiaqi Zhang, Xiaolong Ma, Shiwei Liu, and Lu Yin. Condense, don't just prune: Enhancing efficiency and performance in moe layer pruning. *CoRR*, abs/2412.00069, 2024b.
- Mingyu Cao, Gen Li, Jie Ji, Jiaqi Zhang, Xiaolong Ma, Shiwei Liu, and Lu Yin. Condense, don't just prune: Enhancing efficiency and performance in moe layer pruning, 2025. URL https://arxiv.org/abs/2412.00069.
- I-Chun Chen, Hsu-Shen Liu, Wei-Fang Sun, Chen-Hao Chao, Yen-Chang Hsu, and Chun-Yi Lee. Retraining-free merging of sparse moe via hierarchical clustering, 2025.
- 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. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2024, Bangkok, Thailand, August 11-16, 2024, pp. 1280–1297. Association for Computational Linguistics, 2024a.
- 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, 2024b. URL https://arxiv.org/abs/2401.06066.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang,

Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, and S. S. Li. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. *CoRR*, abs/2501.12948, 2025.

- Zican Dong, Han Peng, Peiyu Liu, Wayne Xin Zhao, Dong Wu, Feng Xiao, and Zhifeng Wang. Domain-specific pruning of large mixture-of-experts models with few-shot demonstrations. *CoRR*, abs/2504.06792, 2025.
- William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *J. Mach. Learn. Res.*, 23:120:1–120:39, 2022.
- Ze-Feng Gao, Peiyu Liu, Wayne Xin Zhao, Zhong-Yi Lu, and Ji-Rong Wen. Parameter-efficient mixture-of-experts architecture for pre-trained language models. In *COLING*, pp. 3263–3273. International Committee on Computational Linguistics, 2022.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.
- Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mixtral of experts. *CoRR*, abs/2401.04088, 2024.
- Sneha Kudugunta, Yanping Huang, Ankur Bapna, Maxim Krikun, Dmitry Lepikhin, Minh-Thang Luong, and Orhan Firat. Beyond distillation: Task-level mixture-of-experts for efficient inference. In *EMNLP* (*Findings*), pp. 3577–3599. Association for Computational Linguistics, 2021.
- Dmitry Lepikhin, HyoukJoong 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. In *ICLR*. OpenReview.net, 2021.
- Mike Lewis, Shruti Bhosale, Tim Dettmers, Naman Goyal, and Luke Zettlemoyer. BASE layers: Simplifying training of large, sparse models. In *ICML*, volume 139 of *Proceedings of Machine Learning Research*, pp. 6265–6274. PMLR, 2021.
- Jia LI, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Costa Huang, Kashif Rasul, Longhui Yu, Albert Jiang, Ziju Shen, Zihan Qin, Bin Dong, Li Zhou, Yann Fleureau, Guillaume Lample, and Stanislas Polu. Numinamath. [https://hf-mirror.com/AI-MO/NuminaMath-1.5] (https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina_dataset.pdf), 2024.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. *arXiv preprint arXiv:2305.20050*, 2023.
- Peiyu Liu, Zikang Liu, Ze-Feng Gao, Dawei Gao, Wayne Xin Zhao, Yaliang Li, Bolin Ding, and Ji-Rong Wen. Do emergent abilities exist in quantized large language models: An empirical study. In *LREC/COLING*, pp. 5174–5190. ELRA and ICCL, 2024.
- Yifei Liu, Li Lyna Zhang, Yi Zhu, Bingcheng Dong, Xudong Zhou, Ning Shang, Fan Yang, and Mao Yang. rstar-coder: Scaling competitive code reasoning with a large-scale verified dataset, 2025. URL https://arxiv.org/abs/2505.21297.

- Xudong Lu, Qi Liu, Yuhui Xu, Aojun Zhou, Siyuan Huang, Bo Zhang, Junchi Yan, and Hongsheng Li. Not all experts are equal: Efficient expert pruning and skipping for mixture-of-experts large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pp. 6159–6172. Association for Computational Linguistics, 2024.
- math ai. Aime25 dataset. https://hf-mirror.com/datasets/math-ai/aime25, 2025.
- open r1. Openr1-math-220k.
 OpenR1-Math-220k, 2025.
 https://hf-mirror.com/datasets/open-r1/
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv e-prints*, 2019.
- Qwen Team. Qwen1.5-moe: Matching 7b model performance with 1/3 activated parameters", February 2024. URL https://qwenlm.github.io/blog/qwen-moe/.
- Qwen Team. Qwen3 technical report, 2025. URL https://arxiv.org/abs/2505.09388.
- Mengzhou Xia, Tianyu Gao, Zhiyuan Zeng, and Danqi Chen. Sheared LLaMA: Accelerating language model pre-training via structured pruning. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=09iOdaeOzp.
- Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy S Liang, Quoc V Le, Tengyu Ma, and Adams Wei Yu. Doremi: Optimizing data mixtures speeds up language model pretraining. *Advances in Neural Information Processing Systems*, 36:69798–69818, 2023.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. A survey of large language models. *CoRR*, abs/2303.18223, 2023.

A LLM USAGE

We use LLMs solely to check and correct grammatical errors in our paper.

B PROOFS

The goal of pruning is to minimize the difference between the output of the pruned model and that of the original model—that is, to minimize the distance $d(X_{L+2}, \hat{X}_{L+2})$. The proof of Theorem 4.1 proceeds in the following three steps.

Proof. We first derive a bound between the original and pruned models when pruning occurs in a single layer. For the i-th layer, where \hat{c} experts are selected from the original c experts, the output discrepancy can be bounded using the Lipschitz constants of the MoE layers from i+1 to L.

$$d((\operatorname{Layer}_{L,c} \circ \cdots \circ \operatorname{Layer}_{i,c} \cdots \circ \operatorname{Layer}_{1,c})(X_1), (\operatorname{Layer}_{L,c} \circ \cdots \circ \operatorname{Layer}_{i,\hat{c}} \cdots \circ \operatorname{Layer}_{1,c})(X_1))$$

$$\leq \operatorname{Lip}(\operatorname{Layer}_{L,c} \circ \cdots \circ \operatorname{Layer}_{i+1,c})d(X_{i+1}, \operatorname{Layer}_{i,\hat{c}}(X_i)).$$
(6)

In the second step, we drive the upper bound of the outputs distance when the number of layers is 3.

$$d((\operatorname{Layer}_{3,c}(X_{2}), (\operatorname{Layer}_{3,c}(\hat{X}_{2}))) = d((\operatorname{Layer}_{3,c} \circ \operatorname{Layer}_{2,c} \circ \operatorname{Layer}_{1,c})(X_{1}), (\operatorname{Layer}_{3,c} \circ \operatorname{Layer}_{2,\hat{c}} \circ \operatorname{Layer}_{1,\hat{c}})(X_{1}))$$

$$\leq d((\operatorname{Layer}_{3,c} \circ \operatorname{Layer}_{2,c} \circ \operatorname{Layer}_{1,c})(X_{1}), (\operatorname{Layer}_{3,c} \circ \operatorname{Layer}_{2,c} \circ \operatorname{Layer}_{1,\hat{c}})(X_{1}))$$

$$+ d((\operatorname{Layer}_{3,c} \circ \operatorname{Layer}_{2,c} \circ \operatorname{Layer}_{1,\hat{c}})(X_{1}), (\operatorname{Layer}_{3,c} \circ \operatorname{Layer}_{2,\hat{c}} \circ \operatorname{Layer}_{1,\hat{c}})(X_{1}))$$

$$\leq \operatorname{Lip}(\operatorname{Layer}_{3,c} \circ \operatorname{Layer}_{2,c})d(\operatorname{Layer}_{1,\hat{c}}(X_{1}), \operatorname{Layer}_{1,c}(X_{1}))$$

$$+ \operatorname{Lip}(\operatorname{Layer}_{3,c})d(\hat{X}_{3}, \operatorname{Layer}_{2,c}(\hat{X}_{2}))$$

$$= \operatorname{Lip}(\operatorname{Layer}_{3,c})d(\hat{X}_{3}, \operatorname{Layer}_{2,c}(\hat{X}_{2}))$$

$$= \sum_{i=2}^{3} \operatorname{Lip}_{i\to 3}d(\hat{X}_{i}, \operatorname{Layer}_{i-1,c}(\hat{X}_{i-1})),$$

$$(7)$$

where the first equality follows from the definitions of X_2 and \hat{X}_2 . The second inequality holds according to the triangle inequality. And the third holds according to the inequality 6. The fourth equality follows from the definitions of $\hat{X}_2 = \operatorname{Layer}_{1,\hat{c}}(X_1)$ and $\hat{X}_1 = X_1$. For convenience, we define $\operatorname{Lip}_{i \to L}$ as the Lipschitz constant from the *i*-th layer to the *L*-th layer, which justifies the final equality.

In the third step, we generalize the result from 3 layers to the case of K+1 layers. Suppose the number of layers K, then there exists the following inequality holds

$$d(\operatorname{Layer}_{K,c}(X_K), \operatorname{Layer}_{K,c}(\hat{X}_K)) \le \sum_{i=2}^K \operatorname{Lip}_{i \to K} d(\hat{X}_i, \operatorname{Layer}_{i-1,c}(\hat{X}_{i-1}))$$
(8)

When the the number of layers is K + 1, we have

$$d(\operatorname{Layer}_{K+1,c}(X_{K+1}), \operatorname{Layer}_{K+1,c}(\hat{X}_{K+1}))$$

$$\leq d(\operatorname{Layer}_{K+1,c} \circ \operatorname{Layer}_{K,c}(X_K), \operatorname{Layer}_{K+1,c} \circ \operatorname{Layer}_{K,c}(\hat{X}_K))$$

$$+ d(\operatorname{Layer}_{K+1,c} \circ \operatorname{Layer}_{K,c}(\hat{X}_K), \operatorname{Layer}_{K+1,c} \circ \operatorname{Layer}_{K,c}(\hat{X}_K))$$

$$\leq \operatorname{Lip}_{K+1} d(\operatorname{Layer}_{K,c}(X_K), \operatorname{Layer}_{K,c}(\hat{X}_K)) + \operatorname{Lip}_{K+1} d(\hat{X}_{K+1}, \operatorname{Layer}_{K,c}(\hat{X}_K))$$

$$= \sum_{i=2}^{K+1} \operatorname{Lip}_{i\to K+1} d(\hat{X}_i, \operatorname{Layer}_{i-1,c}(\hat{X}_i)),$$

$$(9)$$

where the first inequality holds according to the triangle inequality. The second inequality holds according to the Lipschitz constant of (K+1)-th layer. The last inequality follows from inequality 8. Therefore, when the K=L, we have

$$d(X_{L+2}, \hat{X}_{L+2}) = \sum_{i=2}^{L+1} \operatorname{Lip}_{i \to L+1} d(\hat{X}_i, \operatorname{Layer}_{i-1, c}(\hat{X}_{i-1}))$$
(10)