# REAP THE EXPERTS: WHY PRUNING PREVAILS FOR ONE-SHOT MOE COMPRESSION

## **Anonymous authors**

Paper under double-blind review

### **ABSTRACT**

Sparsely-activated Mixture-of-Experts (SMoE) models offer efficient pre-training and low latency but their large parameter counts create significant memory overhead, motivating research into expert compression. Contrary to recent findings favouring expert *merging* on discriminative benchmarks, we demonstrate that expert *pruning* is a superior strategy for generative tasks. We prove that merging introduces an irreducible error by causing a "functional subspace collapse", due to the loss of the router's independent, input-dependent control over experts. Leveraging this insight, we propose Router-weighted Expert Activation Pruning (REAP), a novel pruning criterion that considers both router gate-values and expert activation norms. Across a diverse set of SMoE models ranging from 20B to 1T parameters, REAP consistently outperforms merging and other pruning methods on generative benchmarks, especially at 50% compression. Notably, our method achieves near-lossless compression on code generation and tool-calling tasks with Qwen3-Coder-480B and Kimi-K2, even after pruning 50% of experts.

### 1 Introduction

Interest in the Sparsely-activated Mixture-of-Experts (SMoE) architecture for Large Language Models (LLMs) surged following the release of DeepSeek-V3 (DeepSeek-AI et al., 2024) and other high-quality open-weight SMoE LLMs (Jiang et al., 2024; Meta AI Team, 2025; Yang et al., 2025a; Zeng et al., 2025; Baidu, 2025; Kimi Team et al., 2025). Compared to dense models, the SMoEs offer lower latency and more efficient pre-training (Fedus et al., 2022). However, SMoEs require more parameters than dense models to achieve similar accuracy, resulting in significant memory overhead. Further, expert usage imbalance during inference causes poor accelerator utilization, leading to increased latency or compromises such as dropped tokens (Balmau et al., 2025). Expert usage imbalance also represents an opportunity, motivating prior work which investigates whether experts can be compressed without negatively impairing accuracy (Li et al., 2023; Lu et al., 2024). By eliminating or compressing redundant experts, memory overhead is reduced. A more uniform distribution of expert usage would also improve hardware utilization. Expert compression is particularly valuable for use cases which feature small batch sizes such as local deployments and academic research.

Initial expert compression efforts focused on expert pruning, the removal of experts in their entirety. However, expert pruning is a strong intervention on the model's weights. Techniques such as quantization, low-rank compression, and expert merging also offer memory savings but maintain a lossy representation of the less important experts. Crucially, expert merging has recently been demonstrated to outperform expert pruning when evaluated with perplexity and on Multiple Choice (MC) question answering benchmarks (Li et al., 2023; Liu et al., 2024b). However, an evaluation comparing these methods on generative benchmarks has yet to be conducted. In this work, we demonstrate that — when paired with a suitable saliency criterion — expert pruning outperforms expert merging, particularly on generative benchmark tasks such as code generation, creative writing, and mathematical reasoning. Specifically, our main contributions are as follows:

- We prove that expert merging introduces *irreducible error* due to the loss of the router's independent, input-dependent modulation of the expert outputs resulting in *functional subspace collapse*, substantially reducing the functional output space of the compressed SMoE layer. In contrast, in expert pruned SMoEs the router maintains independent control over the remaining experts;
- We introduce Router-weighted Expert Activation Pruning (REAP), a novel expert pruning saliency criterion, which selects experts to prune which contribute minimally to the layer output by considering both the router gate-values and average activation norm of the experts;
- Across diverse SMoE architectures ranging from 20B to 1T parameters and a suite of generative evaluations, we demonstrate the significant and consistent advantage of REAP over existing expert pruning and

merging approaches, particularly at 50% compression. Notably, our method achieves near-lossless compression on code generation tasks after pruning 50% of experts from Qwen3-Coder-480B and Kimi-K2;

 Our anonymized code is available to facilitate peer-review and we will open source the code and select compressed model checkpoints upon acceptance to facilitate further research on compressed SMoEs and their applications.

# 2 RELATED WORK

**Sparsely activated SMoE architecture.** A Mixture-of-Experts (MoE) layer is comprised of multiple, specialized feed-forward subnetworks known as *experts* and a router which produces gate-values (i.e., *gates*) to dynamically modulate the output of the experts based on the input. The architecture was revived in the deep learning era by the introduction of the SMoE by Shazeer et al. (2017). SMoEs layers only select a subset of experts to use for each input, enabling massive scaling of model parameters without a commensurate increase in computational cost (Lepikhin et al., 2021; Fedus et al., 2022). In transformer-based LLMs, SMoE layers are integrated by replacing the traditional feed-forward layers. Further innovations such as auxiliary-loss-free load balancing (DeepSeek-AI et al., 2024), shared experts, and fined-grained experts (Dai et al., 2024) have propelled SMoE architectures to become the *de facto* standard for LLMs in recent months.

**Expert pruning.** Although SMoE layers effectively decouple total model parameters from inference costs, the memory overhead of storing large SMoEs restricts their deployment in resourced-constrained environments, motivating research in expert pruning to reduce total number of parameters. Early efforts demonstrated that progressively pruning experts based on router weights during fine-tuning until a single expert remained could preserve model quality in task-specific settings (Chen et al., 2022). Koishekenov et al. (2023) found expert pruning to be effective without further fine-tuning despite aggressively pruning up to 80% of experts. Muzio et al. (2024) found that global pruning using gate-values as a saliency criterion was more effective than uniform, layer-wise frequency-based pruning. Other sophisticated pruning criteria have been proposed: Lu et al. (2024) introduced an exhaustive search strategy which prunes experts that minimize the reconstruction loss between the original and pruned layer outputs; Liu et al. (2024a) used a gradient-free evolutionary algorithm to prune experts. Both of these works demonstrated significant improvements over naive frequency-based pruning. A comprehensive evaluation of 16 diverse pruning criteria was conducted by Jaiswal et al. (2025). Expert Activation Norm (EAN) was empirically found to be the highest performing criterion and the benefits of iterative pruning were presented.

Expert merging. While the above-noted works prove that expert compression is feasible via pruning, an alternative compression technique is to *merge* experts. Generally, merging requires both a clustering algorithm and a merging technique. Li et al. (2023) introduced Merge Sparse Mixture of Experts (M-SMoE) which first initializes expert cluster centres by identifying the *dominant* experts with the highest usage frequency globally across all layers. The remaining non-dominant experts are clustered based on the cosine similarity of router logits. Finally, experts weights are aligned via permutation with the weight matching algorithm (Ainsworth et al., 2023) and merged using frequency-weighted parameter averaging. Li et al. (2023) found that their technique outperformed Chen et al.'s (2022) pruning method on MC benchmarks. Chen et al. (2025) proposed Hierarchical Clustering for Sparsely activated Mixture of Experts (HC-SMoE). HC-SMoE clusters experts based on the euclidean similarity of their *representative vectors*—the average activation of each expert measured on *every* token in a calibration dataset — using hierarchical agglomerative clustering. Similar to M-SMoE, HC-SMoE uses frequency-weighted parameter averaging to merge clusters into a single merged expert. Without any fine-tuning, Chen et al. (2025) found that their technique outperformed expert pruning based on router logits (He et al., 2025a), frequency, and Lu et al.'s (2024) method when benchmarked on a suite of MC question answering tasks.

Other compression techniques. In addition to pruning and merging, experts may be compressed through quantization (Huang et al., 2025), low-rank decomposition (Yang et al., 2024a; Gu et al., 2025; He et al., 2025b), weight sparsity (He et al., 2025a), or a combination of any of the above techniques (Liu et al., 2025). These other approaches are orthogonal to expert pruning and merging; however, note that expert merging necessitates re-quantization for block quantization formats that share common scaling coefficients across a group of weights.

**Model merging.** Model merging aims to combine parameters from multiple trained neural networks and has been rapidly adopted as a cost-effective way to improve model quality across diverse domains. The initial motivation for merging was based on the finding that mode connectivity exists between the

loss landscapes of two or more trained neural networks, enabling interpolation of their parameters without incurring an increase in loss (Garipov et al., 2018; Ainsworth et al., 2023; Ito et al., 2024). Simple parameter averaging remains an effective technique; however, more sophisticated strategies based on task vectors have also been proposed to minimize interference in the merged model parameters (Ilharco et al., 2023; Yadav et al., 2023; Yu et al., 2024). Much of the existing literature focuses on the setting in which multiple fine-tunes of a single checkpoint are merged. Non-local merging in which the models do not share a common checkpoint is more closely related to expert merging. Sharma et al. (2024) found that re-scaling of model activations was necessary to achieve high-quality non-local merging.

**LLM evaluation.** Evaluating LLMs is challenging; prior work demonstrated that simple metrics such as perplexity can be misleading when used to evaluate compressed LLMs (Jaiswal et al., 2024). MC benchmarks typically measure the log-likelihood of answer tokens to determine a model's response to a question (Gao et al., 2023; Chandak et al., 2025). As such, each response choice is evaluated in a single forward pass, without any tokens being generated by the model. Perplexity and MC accuracy can therefore be viewed as discriminative metrics. In contrast, generative benchmarks require the model to output a response, more closely corresponding with real-world use-cases of LLMs. Tasks such as code generation, mathematical reasoning with structured outputs, and creative writing are examples of generative benchmarks.

# MERGING EXPERTS CAUSES FUNCTIONAL SUBSPACE COLLAPSE

**Setup.** To motivate our proposed expert pruning method, we first formally develop the expected errors of both expert merging and pruning. Consider a SMoE layer with K experts  $f_1,...,f_K$ , each a function  $f_k:\mathbb{R}^d\to\mathbb{R}^d$ , and a router producing non-negative gates  $g(x)=(g_1(x),...,g_K(x))\in\Delta^{K-1}$ . Top-k routing is achieved by zeroing all but the largest k gates. The output of the original layer is

$$h(x) := \sum_{k=1}^{K} g_k(x) f_k(x). \tag{1}$$

**Two operations at fixed compression.** To analyse the fundamental difference between compression operations, we focus on the elementary case of reducing two experts,  $(f_i, f_j)$ , to one. This pairwise analysis is the building block for any larger merge within a cluster. Pruning removes expert j and re-normalizes the router outputs over the remaining K-1 experts, producing a new set of gates  $\bar{g}(x)$ . Merging replaces  $(f_i, f_j)$  with a new expert f. Existing one-shot expert merging methods such as HC-SMoE and M-SMoE sum the gates for the original experts  $g_i(x) + g_j(x)$ . The pruned,  $\bar{h}(x)$ , and merged,  $\bar{h}(x)$ , layer outputs are

$$\bar{h}(x) := \sum_{k \neq j} \bar{g}_k(x) f_k(x),$$
 (2)  $\tilde{h}(x) := \sum_{k \neq i,j} g_k(x) f_k(x) + (g_i(x) + g_j(x)) \tilde{f}(x).$  (3)

# 3.1 MERGING INDUCES AN INPUT-DEPENDENT TARGET A SINGLE EXPERT CANNOT REALIZE

Define the router's input-dependent mixing ratio  $r(x) := \frac{g_i(x)}{g_i(x) + g_j(x)} \in [0,1]$  on the set where  $g_i + g_j > 0$ . Substituting  $g_i(x)$  and  $g_j(x)$  in terms of r(x), the original contribution of the pair (i,j) can be written as

$$g_{i}(x)f_{i}(x) + g_{j}(x)f_{j}(x) = \left[r(x)(g_{i}(x) + g_{j}(x))\right]f_{i}(x) + \left[(1 - r(x))(g_{i}(x) + g_{j}(x))\right]f_{j}(x)$$

$$= \left(g_{i}(x) + g_{j}(x)\right)\underbrace{\left(r(x)f_{i}(x) + \left(1 - r(x)\right)f_{j}(x)\right)}_{\text{The ideal, input-dependent target expert}}.$$

$$\tag{4}$$

After merging, the router must apply the summed gate,  $g_i(x) + g_j(x)$ , to a constant convex combination of the constituent experts which is independent of x. The core issue is that the merged model is forced to approximate the dynamic, input-dependent target expert with a static one. The following theorem quantifies this unavoidable approximation error.

**Theorem 1** (Irreducible error of merging). Let  $\hat{f}_{\alpha}(x) = \alpha f_i(x) + (1-\alpha) f_i(x)$  with a constant  $\alpha \in [0,1]$ and define  $\Delta_{ij} := f_i(x) - f_j(x)$ . The  $L^2$  error of the merged pair is minimized when  $\alpha$  is chosen to be

the expected mixing ratio,  $\alpha^* := \mathbb{E}[r(x)]$ . Omitting the argument (x) for brevity, this minimal error is

$$\|(g_i+g_j)(rf_i+(1-r)f_j)-(g_i+g_j)(\alpha f_i+(1-\alpha)f_j)\|^2 = \underbrace{\mathbb{E}[(g_i+g_j)^2]}_{router \ scale} \cdot \underbrace{\operatorname{Var}[r(x)]}_{policy \ variability} \cdot \underbrace{\|\Delta_{ij}\|^2}_{expert \ gap}. \tag{5}$$

In particular, if the router's policy is not constant  $(\operatorname{Var}[r(x)] > 0)$  and the experts are not functionally identical  $(\|\Delta_{ij}\| > 0)$ , then every constant- $\alpha$  merge incurs strictly positive excess risk.

*Proof.* The error term simplifies to  $\|(g_i+g_j)(r-\alpha)\Delta_{ij}\|^2$ . Assuming independence between the router policy and expert functions, this is proportional to  $\mathbb{E}[(r-\alpha)^2]$ . This is a standard least-squares problem minimized when  $\alpha = \mathbb{E}[r]$ , and the minimal value is  $\operatorname{Var}[r]$ .

**Consequences.** Theorem 1 illustrates that merging with summed gates is fundamentally flawed whenever (i) the router has learned an input-dependent policy for mixing two experts (Var[r] > 0), and (ii) the experts are themselves distinct  $(\|\Delta_{ij}\| > 0)$ . Any fixed  $\alpha$  cannot overcome the irreducible error bound established in equation 5.

#### 3.2 Pruning preserves independent control

Pruning removes one function but importantly does *not* tie the remaining gates. The router still modulates each surviving expert *independently*. In contrast, merging removes a degree of freedom in the policy by replacing individual experts with their mergers. For a direct comparison under no fine-tuning, pruning expert *j* and reallocating its gate-value to expert *i* produces the error

$$\|(g_i(x)f_i(x)+g_j(x)f_j(x))-(g_i(x)+g_j(x))f_i(x)\|^2 = \mathbb{E}[g_j(x)^2\|\Delta_{ij}(x)\|_2^2].$$
 (6)

Unlike equation 5, equation 6 *does not* penalize policy variability, the router still controls surviving experts independently. Whenever the router substantially mixes i and j (large Var[r]) while the pruned expert j has a small average gate-value ( $\mathbb{E}[g_i^2]$ ), pruning admits a strictly smaller error than merging.

**Synthesis.** Theorem 1 establishes that summed gate merging incurs an irreducible error directly proportional to the router's policy variability  $(\operatorname{Var}[r(x)])$ . In contrast, the error from pruning a low-usage expert (Eq. 6) is proportional to its gate-value  $(\mathbb{E}[g_j^2])$  and is insensitive to policy variability. Therefore, when the router actively mixes between two distinct experts, merging is fundamentally disadvantaged.

**Remarks.** (i) The constant-mixture model  $\tilde{f}_{\alpha}$  is mathematically related to the frequency weighted parameter averaging merge used in practice. (ii) Even if  $\tilde{f}$  was dependent on x, the router after merging cannot independently modulate the two latent functions, so the original policy is invalidated. (iii) With top-k routers, the specific irreducible error from policy variability  $(\mathrm{Var}[r(x)])$  is generated exclusively on the support where both experts are selected. Outside that support, this component vanishes, leaving only a static error term that depends on the functional expert gap. (iv) See Appendix A for an extension of the above analysis to hierarchical clustering.

## 3.3 EMPIRICAL EVIDENCE FOR LOSS OF INDEPENDENT CONTROL

**Setup.** We analyse the functional subspaces of expert outputs across four diverse state-of-the-art SMoE architectures by recording mean expert activations from 32 samples of 2048 tokens from the c4 dataset (Raffel et al., 2020). By projecting expert activations onto their first two principal components, we visualize how pruning and merging affect the learned representations. See Appendix B for additional discussion.

**Early vs. late behaviour.** Figures 1 and A4 demonstrate a striking progression of functional collapse from early to late layers across all architectures. In early layers, the original experts form relatively compact manifolds with moderate spread. After pruning, the surviving experts maintain their positions on the original manifold, preserving its geometric structure with reduced density. In contrast, merging produces a visible contraction toward the manifold's centre. The contrast becomes dramatic in late layers, where experts are more specialized, and in high granularity architectures with many experts per layer.

The progression from early to late layers validates our theoretical prediction that the irreducible error is proportional to Var[r(x)]. Early layers, which typically learn more generic features, exhibit lower

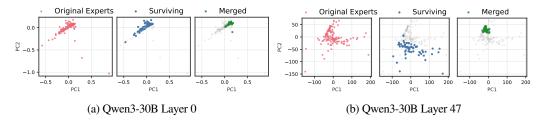


Figure 1: (a) Functional subspace (PCA) for early SMoE layers in Qwen3-30B. Pruning (blue) preserves the manifold geometry; merging (green) collapses it toward the centre. (b) Functional subspace (PCA) for late MoE layers. The contraction under merging is dramatically more pronounced, with up to  $100 \times$  reduction in spread for models with many experts. See Figure A4 for results from other models.

policy variability and thus less dramatic collapse. Late layers, where experts have specialized for distinct computational roles, demonstrate high policy variability, resulting in the severe functional collapse observed when these specialized experts are merged into static averages.

**Synthesis across architectures.** The consistency of these patterns across architectures with vastly different expert counts (8 to 128), sparsity levels (6.25% to 25% active), and parameter scales (21.9B to 109B) demonstrates that functional collapse under merging is a fundamental property of the operation rather than an artifact of specific implementations. These visualizations reveal that the core issue is not the reduction in the number of experts *per se*, but rather the qualitative change in the router's control structure.

# 4 ROUTER-WEIGHTED EXPERT ACTIVATION PRUNING (REAP)

The above analysis demonstrates that the functional output space of a SMoE layer is defined by the *coordinated behaviour* of the router and experts. An expert's total contribution to its layer's output is determined by both its gate-value,  $g_k(x)$ , and the magnitude of its output vector,  $\|f_k(x)\|_2$ . However, naive frequency-based pruning fails to consider these properties. Intuitively, pruning experts which contribute minimally to the layer output minimizes the difference between the original and pruned layer outputs. Let h(x) be the original output and  $\bar{h}_{\backslash j}(x)$  be the output after pruning expert j and re-normalizing the remaining router weights. The error induced by pruning expert j is

$$\Delta \bar{h}_{\backslash j}(x) := h(x) - \bar{h}_{\backslash j}(x) = \sum_{k} g_k(x) f_k(x) - \sum_{k \neq j} \frac{g_k(x)}{1 - g_j(x)} f_k(x). \tag{7}$$

Re-normalization of the router weights after pruning expert j modulates all other remaining expert outputs, making direct minimization of  $\Delta h_j$  complex. However, since our goal is to prune unimportant experts, we can reasonably assume their gate-values are small when active  $\mathbb{E}_{x\sim\mathcal{X}}[g_j(x)]\ll 1$ . Under this assumption, the weight re-normalization factor is negligible, i.e.,  $1-g_j(x)\approx 1$ , and the error induced by pruning expert j is approximately equal to the expert's direct contribution to the layer output

$$\Delta \bar{h}_{\backslash j}(x) \approx \sum_{k} g_k(x) f_k(x) - \sum_{k \neq j} g_k(x) f_k(x) = g_j(x) f_j(x). \tag{8}$$

To select which experts to prune, we propose a novel saliency criterion, REAP, which approximates an expert's importance by measuring its direct contribution to the layer's output magnitude. Specifically, the saliency score,  $S_i$ , is defined as the average of this contribution over tokens for which the expert is active

$$S_j = \frac{1}{|\mathcal{X}_j|} \sum_{x \in \mathcal{X}_j} g_j(x) \cdot \left\| f_j(x) \right\|_2, \tag{9}$$

where  $\mathcal{X}_j$  is the set of tokens where expert j is active (i.e.,  $\mathcal{X}_j = \{x \mid j \in \text{TopK}(\mathbf{g}(x))\}$ ). The experts with the minimum saliency score are selected for pruning. REAP is robust to outlier activations and has a direct, intuitive interpretation by quantifying the average magnitude an expert adds to the output vector when it is selected by the router. Pruning experts with the lowest  $S_j$  removes those with the least impactful contribution.

270 271

Table 1: Comparison of SMoE models included in our study.

Parameters

Active

First layer

Shared

Routed

277278279280

280281282

283

284

285286287288289

290

291

292293294295296

297

298

299

300

306

307

308

313

314

315

316

317

318 319

320

321

322

323

Model Top-K Sparsity Params. (1e9) Experts Experts (1e9)dense ERNIE-4.5-21B-A3B-PT 87.88% 21.9 64 2 3 Yes Owen3-30B-A3B 128 0 93.75% 30.5 3 8 No Mixtral-8x7B-Instruct-v0.1 8 0 2 75.00% 46.7 13 No GLM-4.5-Air 128 1 8 93.02% 106.9 12 Yes Llama-4-Scout-17B-16E-Instruct 16 1 1 88.24% 107.8 17 No Qwen3-Coder-480B-A35B-Instruct-FP8 160 0 8 95.00% 480.2 35 No Kimi-K2-Instruct-W4A16 (RedHatAI, 2025) 97.66% 1026.4 32 384 1 8 Yes

## 5 EXPERIMENTS

**Setup.** We implement REAP and other expert compression baselines in PyTorch (Ansel et al., 2024). We collect router logits and expert activation data to calibrate the compression algorithms using a variety of general pre-training and domain-specific Supervised Fine-Tuning (SFT) datasets. For calibration, 1,024 samples are randomly selected and packed to 2,048 sequence length for models with  $\leq$  110B parameters. For models with  $\geq$  110B parameters, we select 12,228 samples with a maximum sequence length of 16,384 tokens without truncation or packing.

We compress models by pruning or merging 25% or 50% of experts in each layer, except for M-SMoE which determines the number of clusters per layer based on global expert usage frequency. When evaluating models with  $\leq 50$ B parameters on coding and MC, we calibrate and compress the models using three different seeds and report the mean. Larger models, creative writing, and mathematical reasoning evaluations are reported using a single seed, except where explicitly noted otherwise. All models are evaluated in the one-shot setting, with no additional fine-tuning after compression.

Models and data. We evaluate the expert compression algorithms on a diverse set of six SMoE architectures covering model sizes from 21B to 1T with varying degrees of sparsity and expert granularity, see Table 1 for details. For MC question answering and code generation benchmarks, we use c4 (Raffel et al., 2020; Allen Institute for AI, 2024) and evol-codealpaca (Chaudhary, 2023; Luo et al., 2024; Tam, 2023) datasets to assess both general and domain-specific calibration. Models with ≥ 110B parameters are additionally calibrated with data from xlam-function-calling (Liu et al., 2024c; Salesforce, 2025) and SWE-smith-trajectories (Yang et al., 2025c;b) datasets. For creative writing and math benchmarks we employ WritingPrompts curated (Pritsker, 2024) and tulu-3-sft-personas-math (Lambert et al., 2025; Allen Institute for AI, 2025), respectively. The default chat template is applied to all SFT datasets and

**Evaluation.** Compressed SMoE models are evaluated on a suite of benchmarks including MC question answering, code generation, mathematical reasoning, creative writing, and tool calling. See Appendix C for details. We implement HC-SMoE and M-SMoE as expert merging baselines. Average linkage criterion is used for HC-SMoE. M-SMoE does not include low-rank compression from the complete MC-SMoE method. Pruning baselines consist of frequency-based pruning and EAN. See Appendix D for formal definitions.

## 5.1 RESULTS

In Table 2 and Figure 2 code generation, creative writing, math reasoning, and MC results are presented for Qwen3-30B and GLM-4.5-Air after calibration with the evol-codealpaca dataset. Table 3 contains results for large-scale SMoE pruned models on code generation and tool calling benchmarks. See Table A4 and Table A5 for detailed MC and code generation results, respectively. Figure A5 depicts coding generation and MC accuracy verses model parameters. See Appendix E for additional results.

**Zero-shot MC question answering.** Both merging and pruning are capable of producing accurate compressed SMoE models for MC question answering. HC-SMoE and REAP have a mean decrease in accuracy of approximately 4% and 13% for compression ratios of 25% and 50%, respectively, excluding large-scale SMoEs. REAP achieves first or second rank among all methods, models and compression ratios, suggesting strong consistency regardless of specific model architecture. When calibrated on c4, we find slightly improved accuracies for all compression methods with similar rankings as noted above, see Table A6.

324 325 326

Table 2: MC and generative benchmark results for Qwen3-30B and GLM-4.5-Air.

341342343344

345

340

352353354

355356357358

359

360

361

362

363

364

365366367368369

370

371

372

374

375

376

377

Coding Creative Writing WildBench Math MC Model GSM8K MATH-500 Math Avg MC Avg Compression Technique Method Eval-LiveCode Code Avg 0.872 0.721 Baseline 0.859 0.302 0.581 0.811 0.903 0.887 M-SMoE 0.822 0.293 0.558 0.805 0.901 0.872 0.886 0.558 Merging HC-SMoE 0.258 0.674 0.800 0.497 25% Frequency 0.302 0.576 0.905 0.864 0.600 0.849 0.807 0.885 0.811 0.804 0.881 0.878 0.311 0.576 0.866 0.864 Owen3-30B-A3B REAP 0.843 0.892 0.575 0.669 0.413 0.379 0.451 M-SMoF 0.62 0.205 0.725 0.824 0.838 0.831 Merging 0.574 HC-SMoE 0.542 0.185 0.008 0.760 0.696 0.728 50% 0.236 0.470 0.677 0.882 0.860 0.483 0.704 Pruning 0.702 0.493 0.798 0.306 0.886 0.842 0.552 **0.557** REAP 0.821 0.293 0.878 0.872 0.875 0.374 0.597 0.839 0.846 0.918 0.882 Baseline 0.820 0.747 M-SMoE HC-SMoE 0.781 0.793 0.330 0.555 0.578 0.781 0.788 0.848 0.842 Merging 0.704 0.908 0.875 25% Frequency 0.341 0.793 0.832 0.908 0.597 Pruning 0.374 GI M-4 5-Air 0.821 0.824 0.839 0.908 0.874 0.637 0.880 REAP 0.390 0.831 0.592 0.835 0.678 M-SMoF 0.493 0.099 0.296 0.391 0.465 0.466 0.465 0.444 Merging HC-SMoE 0.441 0.593 0.700 0.564 50% 0.104 0.325 0.604 0.615 0.612 0.613 0.521 Frequency 0.546 Pruning 0.838 0.511 0.513 0.553 0.255 REAP 0.559

Table 3: Benchmark results for agentic, non-agentic coding and tool-use tasks.

Model Compression		Method	No Eval+	on-Agentic C LiveCode	Coding Code Avg	Agentic Coding SWE-Bench-Verified	Non-Live	Tool Us Live	e (BFCLv3) Multi-Turn	Overall
Qwen3-Coder- 480B-A35B- Instruct-FP8	Baseline		0.889	0.431	0.660	0.540	0.866	0.825	0.380	0.690
	25%	Frequency EAN REAP	0.792 0.876 0.884	0.296 0.419 0.416	0.544 0.647 <b>0.650</b>	0.378 0.534 <b>0.540</b>	0.844 0.831 0.878	0.763 0.813 0.823	0.355 0.384 0.392	0.654 0.676 <b>0.698</b>
	50%	EAN 0.831 0.38		0.012 0.382 0.415	0.011 <u>0.607</u> <b>0.644</b>	0.000 <b>0.536</b> <u>0.522</u>	0.200 0.822 0.849	0.392 0.774 0.801	0.000 0.383 0.371	0.197 0.659 <b>0.674</b>
	Baseline		0.883	0.434	0.659	0.554	0.840	0.802	0.355	0.666
Kimi-K2- Instruct- W4A16	25%	Frequency EAN REAP	0.524 0.831 0.889	0.082 0.379 0.440	0.303 0.605 <b>0.664</b>	0.000 <u>0.562</u> <b>0.580</b>	0.644 0.819 0.842	0.603 0.802 0.801	0.045 0.335 0.263	0.431 <b>0.652</b> <u>0.635</u>
WHAIO	50%	Frequency EAN REAP	0.124 0.772 0.863	0.000 0.253 0.429	0.062 0.513 <b>0.646</b>	0.000 0.576 0.576	0.255 0.778 0.785	0.397 0.767 0.743	0.003 0.173 0.164	0.218 <b>0.573</b> <u>0.564</u>

**Generative benchmarks.** Compared to MC, generative benchmarks are more representative of real-world use cases of LLMs. In this setting, pruning emerges as the clearly superior compression method on the generative task benchmarks. Excluding large-scale SMoEs, REAP achieves a mean decrease in accuracy of 2.8% and 8.0% at 25% and 50% compression ratios, respectively, on coding. In comparison, both HC-SMoE and M-SMoE produce mean decreases in accuracy >5% at 25% compression and >20% at 50% compression. Notably, REAP maintains significantly higher accuracy at 50% compression than other pruning methods. On creative writing, REAP and EAN are near-lossless at 25% compression with REAP offering improved quality at 50% compression. Merging methods are less consistent across various model architectures and compression ratios. For example, M-SMoE is the best method for Qwen3-30B at 50% compression, but the worst on GLM-4.5-Air. REAP attains the best mathematical reasoning results with a remarkable mean decrease in accuracy of just 1.1% at 50% compression. HC-SMoE and M-SMoE offer high accuracy at 25% compression but are significantly less accurate than pruning at 50% compression.

**Expert pruning at scale.** To asses whether pruning remains viable at scale, we prune Qwen3-Coder-480B and Kimi-K2-Instruct. On MC questions, REAP outperforms other pruning methods. On non-agentic coding tasks, REAP achieves near-lossless accuracy with a 0.20% and 1.4% mean decrease in accuracy compared to baseline at 25% and 50%, respectively, outperforming EAN and frequency-based pruning, particularly at 50% compression. On the challenging SWE-Bench task, both REAP and EAN maintain high accuracy at 25% and 50% compression, with some scores slightly exceeding the baseline. On tool use, EAN and REAP are comparable, with REAP slightly outperforming at 50% compression with a mean decrease in accuracy of 5% versus 6% for EAN. Frequency-based pruning suffers from a sharp degradation in quality at 50% compression, highlighting the importance of pruning saliency criteria which consider expert activations. Scaling the pruning methods is relatively trivial. Unlike HC-SMoE, calibration for pruning does not require recording activations from every expert for every token, facilitating efficient

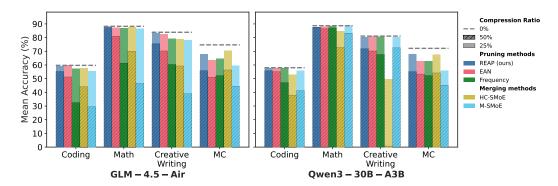


Figure 2: GLM-4.5-Air and Qwen3-30B accuracy vs. task type. REAP offers significant improvements compared to other methods at 50% compression. Note the significant performance drop for merging methods on generative tasks (Coding, Math, Creative Writing) compared to their relative strength on MC.

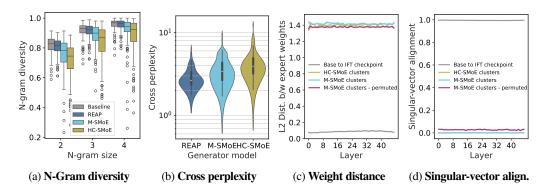


Figure 3: (a) & (b) **N-Gram diversity** and **cross-perplexity** of compressed Qwen3-30B-A3B models at 50% compression, respectively. (c) & (d) The mean relative **L2-distance** and **singular-vector alignment** between expert weights Qwen3-30B at 50% compression, respectively.

calibration. Further, pruning can be easily applied to quantized models without any additional steps required to reconcile block scales or re-quantize following compression.

Quantifying merged MoE generation quality. While merged expert SMoEs offer reasonable quality for discriminative tasks such as MC question and answering, they fail to remain competitive on generative tasks. To help explain the performance gap of merged models between discriminative and generative tasks, we perform an analysis of the compressed model outputs and compare with REAP pruned models. We prompt 50% compressed Qwen3-30B models with 100 questions randomly sampled from the evol-codealpaca dataset and record their outputs. In Figure 3a, we measure the N-gram diversity and find that the merged models have significantly lower diversity across all N-gram sizes measured. In contrast, the REAP pruned model remains similar to the base model, albeit slightly less diverse. In Figure 3b, we measure the perplexity of the text generated by the compressed models with the original baseline model. The text generated by the merged models has both a higher mean and higher variance than the pruned model generations, suggesting that the REAP pruned model outputs are more closely aligned to the original model.

**The challenges of expert merging.** Model merging has been widely adopted to facilitate LLM fine-tuning. Why does expert merging miss the mark? In addition to the loss of the router's input-dependent modulation of experts explored in Section 3, we argue that the non-local nature of expert merging and high cardinality of expert clusters pose significant unresolved challenges.

In Figure 3c, we plot the mean relative L2-distance between experts clustered by HC-SMoE or M-SMoE and compare with the distance between expert weights from the pretrained to Instruct Fine-Tuned (IFT) checkpoints. We find that the distance between clustered experts within the same layer greatly exceeds that of experts in the IFT checkpoint after fine-tuning. Ito et al. (2024) found that weight matching permutations improved alignment of parameters' singular vectors. Following their approach, we decompose expert

weights with Singular Value Decomposition (SVD) and plot the singular-vector alignment in Figure 3d. Even after applying weight matching permutations, the M-SMoE expert clusters remain far apart both in weight space and singular-vector alignment. The relatively poorly aligned experts highlight the considerable challenge of coherently merging their parameters.

When merging works well, it's more closely related to pruning than one might expect. In Figure A6a, we depict the frequency of singleton clusters — clusters containing a single expert — for both HC-SMoE and M-SMoE. A singleton cluster is directly analogous to an expert that remains after pruning. We find that HC-SMoE in particular has a high prevalence of singleton clusters, leaving important experts unadulterated and compressing the rest into a few *mega*-clusters containing tens of experts. This is particularly true of the high granularity models which contain more experts per layer. We hypothesize that the cardinality of these mega-clusters poses a challenge for existing merging algorithms and test this intuition in Figure A6b. Unfortunately, even modest restrictions of the maximum cluster size to 32 — half the number of experts to compress — results in large decreases in model quality on coding tasks.

The importance of domain-specific calibration. In Figure A7, we plot the code generation accuracy of the various compression methods and models when calibrated on either c4 or evol-codealpaca. The difference is stark, c4 calibration results in a collapse in accuracy, with several compressed model instances failing to produce coherent outputs, resulting in 0% accuracy. In Figure A8, we compare the accuracy of compressed Qwen3-30B models calibrated with either domain-specific data or the combined calibration data across all generative tasks. The domain-specific calibrated models achieve significantly higher accuracy, especially at 50% compression. At high compression ratios, domain-specific calibration is crucial.

# 6 DISCUSSION

Similar to prior work, we find that expert merging performs reasonably well on MC benchmarks. This may be because MC tasks only require a discriminative function that can be approximated by an *average* expert. In contrast, merging fails to maintain model quality on generative tasks, particularly at 50% compression. Generative tasks require auto-regressive generation, a capability that is lost when the router's fine-grained control is removed. Compared to expert pruning, merging is less consistent, exhibiting higher variance across models and compression ratios. The outputs of expert merged models are more repetitive and less closely aligned with the base model compared with pruned models. Taken together, these observations are direct evidence of alterations to the functional manifold of the SMoE layers discussed in Section 3.3 stemming from the loss of the router's input-dependent control over the experts and subsequent introduction of novel functions due to tying of the merged expert gates.

Overall, expert pruned models offer consistently higher accuracy than merged models on generative tasks. REAP is a robust pruning criterion that generalizes across a wide array of SMoE architectures, compression ratios, and generative tasks. By taking into consideration both the router gate-values and expert activation norms, REAP prunes the experts which contribute the least to each layers output on a per-token average, regardless of usage frequency. REAP is scalable, achieving near-lossless compression on coding tasks with Qwen3-Coder-480B and Kimi-K2. The successes of REAP highlights the crucial importance of preserving coordination between the router and experts. Compression methods which impair the router's ability to independently modulate expert outputs or tie gate-values are less likely to succeed.

Finally, this work highlights the importance of comprehensive downstream evaluations and the significant challenges involved with evaluating LLMs. Discriminative metrics such as perplexity and log-likelihood based MC benchmarks are not necessarily good proxies for generative model quality.

# 7 CONCLUSION

Our analysis of current SMoE expert merging techniques finds that the router's loss of independent control over experts results in *functional subspace collapse*. In contrast, expert pruning produces a coordinate subspace of the original layer which maintains the topology of the functional manifold. Based on our findings that the coordination between the router and experts is fundamental, we introduce REAP, a novel expert pruning method which prunes experts that contribute the least to the layer's output. Empirically, we demonstrate that REAP retains remarkably high accuracy on an wide array of generative tasks across a diverse set of model architectures. We hope that this work inspires further compression techniques for SMoEs and facilitates the deployment of accurate, domain-specific models in resource constrained settings.

### **ETHICS STATEMENT**

This work research focused on the algorithmic compression of SMoE models and does not involve the use of human subjects, personally identifiable information, or sensitive data. The datasets used for calibration and evaluation (e.g., c4, evol-codealpaca) are publicly available. Our aim is enable the use of large-scale SMoE models in resource constrained settings. However, we acknowledge that compression techniques such as REAP could potentially facilitate deployment of models for malicious purposes. Further, our compression methods are applied to pre-trained models and any biases related to fairness, discrimination, or representation inherent in the original models may be present in their compressed versions. We make no attempt in this work to mitigate these potential biases. The primary contribution of this paper is technical, and we do not foresee any new, direct ethical concerns arising from our proposed methodology beyond those already associated with the deployment of large language models.

### REPRODUCIBILITY STATEMENT

We are committed to ensuring the reproducibility of our research. As stated in the introduction, we have submitted anonymized code to facilitate the peer-review process and will release the full source code and select compressed model checkpoints upon the paper's acceptance. REAP is formally described in Section 4. The baseline methods we compare against, including frequency-based pruning, EAN, M-SMoE, and HC-SMoE, are formally defined in Appendix D. Section 5 provides a detailed description of our experimental setup, including the specific models used, the calibration and evaluation datasets, and the implementation details for all compression experiments. Further evaluation details are provided in Appendix C.

#### REFERENCES

Samuel Ainsworth, Jonathan Hayase, and Siddhartha Srinivasa. Git re-basin: Merging models modulo permutation symmetries. In *Proceedings of the Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=CQsmMYmlP5T.

Allen Institute for AI. allenai/c4·Datasets at Hugging Face, August 2024. URL https://huggingface.co/datasets/allenai/c4.

Allen Institute for AI. allenai/tulu-3-sft-personas-math · Datasets at Hugging Face, 2025. URL https://huggingface.co/datasets/allenai/tulu-3-sft-personas-math.

Jason Ansel, Edward Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael Voznesensky, Bin Bao, Peter Bell, David Berard, Evgeni Burovski, Geeta Chauhan, Anjali Chourdia, Will Constable, Alban Desmaison, Zachary DeVito, Elias Ellison, Will Feng, Jiong Gong, Michael Gschwind, Brian Hirsh, Sherlock Huang, Kshiteej Kalambarkar, Laurent Kirsch, Michael Lazos, Mario Lezcano, Yanbo Liang, Jason Liang, Yinghai Lu, CK Luk, Bert Maher, Yunjie Pan, Christian Puhrsch, Matthias Reso, Mark Saroufim, Marcos Yukio Siraichi, Helen Suk, Michael Suo, Phil Tillet, Eikan Wang, Xiaodong Wang, William Wen, Shunting Zhang, Xu Zhao, Keren Zhou, Richard Zou, Ajit Mathews, Gregory Chanan, Peng Wu, and Soumith Chintala. PyTorch 2: Faster Machine Learning Through Dynamic Python Bytecode Transformation and Graph Compilation. In 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2 (ASPLOS '24). ACM, April 2024. doi: 10. 1145/3620665.3640366. URL https://docs.pytorch.org/assets/pytorch2-2.pdf.

Baidu. Ernie 4.5 technical report, 2025. URL https://yiyan.baidu.com/blog/publication/ERNIE\_Technical\_Report.pdf.

Oana Balmau, Anne-Marie Kermarrec, Rafael Pires, André Loureiro Espírito Santo, Martijn de Vos, and Milos Vujasinovic. Accelerating moe model inference with expert sharding. In *Proceedings of the 5th Workshop on Machine Learning and Systems*, pp. 192–199, 2025.

Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo. The fifth pascal recognizing textual entailment challenge. *TAC*, 7(8):1, 2009.

Nikhil Chandak, Shashwat Goel, Ameya Prabhu, Moritz Hardt, and Jonas Geiping. Answer Matching Outperforms Multiple Choice for Language Model Evaluation, July 2025. URL http://arxiv.org/abs/2507.02856. arXiv:2507.02856 [cs].

541

542

543

544

546 547

548

549

550

551

552

553

554

555

556

558

559

561

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

588

589

592

Sahil Chaudhary. Code alpaca: An instruction-following llama model for code generation. https://github.com/sahil280114/codealpaca, 2023.

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. In *Proceedings of the Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=hslozRxzXL.

Tianyu Chen, Shaohan Huang, Yuan Xie, Binxing Jiao, Daxin Jiang, Haoyi Zhou, Jianxin Li, and Furu Wei. Task-Specific Expert Pruning for Sparse Mixture-of-Experts, June 2022. URL http://arxiv.org/abs/2206.00277. arXiv:2206.00277 [cs].

Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. BoolQ: Exploring the surprising difficulty of natural yes/no questions. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 2924–2936, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1300. URL https://aclanthology.org/N19-1300/.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge, 2018. URL https://arxiv.org/abs/1803.05457. arXiv:1803.05457 [cs].

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021. URL https://arxiv.org/abs/2110.14168.arXiv:2110.14168 [cs].

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, January 2024. URL http://arxiv.org/abs/2401.06066. arXiv:2401.06066 [cs] version: 1.

DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, 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, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junije Oiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanjia Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. DeepSeek-V3 Technical Report, December 2024. URL http://arxiv.org/abs/2412.19437. arXiv:2412.19437 [cs] version: 1.

- William Fedus, Barret Zoph, and Noam Shazeer. Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity, June 2022. URL http://arxiv.org/abs/2101.03961.arXiv:2101.03961 [cs].
- Leo Gao. Multiple Choice Normalization in LM Evaluation, October 2021. URL https://blog.eleuther.ai/multiple-choice-normalization/.
- 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. A framework for few-shot language model evaluation, 12 2023. URL https://github.com/EleutherAI/lm-evaluation-harness.
- Timur Garipov, Pavel Izmailov, Dmitrii Podoprikhin, Dmitry P Vetrov, and Andrew G Wilson. Loss surfaces, mode connectivity, and fast ensembling of dnns. *Advances in neural information processing systems*, 31, 2018.
- Hao Gu, Wei Li, Lujun Li, Zhu Qiyuan, Mark G. Lee, Shengjie Sun, Wei Xue, and Yike Guo. Delta decompression for moe-based LLMs compression. In *Proceedings of the Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=ziezViPoN1.
- Shwai He, Daize Dong, Liang Ding, and Ang Li. Towards Efficient Mixture of Experts: A Holistic Study of Compression Techniques, March 2025a. URL http://arxiv.org/abs/2406.02500. arXiv:2406.02500 [cs] version: 3.
- Yifei He, Yang Liu, Chen Liang, and Hany Hassan Awadalla. Efficiently Editing Mixture-of-Experts Models with Compressed Experts, March 2025b. URL http://arxiv.org/abs/2503.00634.arXiv:2503.00634 [cs].
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *Proceedings of the Ninth International Conference on Learning Representations*, 2021a. URL https://openreview.net/forum?id=d7KBjmI3GmQ.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In *Proceedings of the Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021b. URL https://openreview.net/forum?id=7Bywt2mQsCe.
- Wei Huang, Yue Liao, Jianhui Liu, Ruifei He, Haoru Tan, Shiming Zhang, Hongsheng Li, Si Liu, and XIAOJUAN QI. Mixture compressor for mixture-of-experts LLMs gains more. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=hheFYjOsWO.
- Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Ludwig Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=6t0Kwf8-jrj.
- Akira Ito, Masanori Yamada, and Atsutoshi Kumagai. Analysis of Linear Mode Connectivity via Permutation-Based Weight Matching: With Insights into Other Permutation Search Methods. In *Proceedings of the Forty-Second International Conference on Machine Learning*, October 2024. URL https://openreview.net/forum?id=lYRkGZZi9D.
- 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. In *Proceedings of the Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=chfJJYC3iL.
- Ajay Jaiswal, Jianyu Wang, Yixiao Li, Pingzhi Li, Tianlong Chen, Zhangyang Wang, Chong Wang, Ruoming Pang, and Xianzhi Du. Finding Fantastic Experts in MoEs: A Unified Study for Expert Dropping Strategies and Observations, April 2025. URL http://arxiv.org/abs/2504.05586.arXiv:2504.05586 [cs].

Ajay Kumar Jaiswal, Zhe Gan, Xianzhi Du, Bowen Zhang, Zhangyang Wang, and Yinfei Yang. Compressing llms: The truth is rarely pure and never simple. In *Proceedings of the Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=B9klVS7Ddk.

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, January 2024. URL http://arxiv.org/abs/2401.04088. arXiv:2401.04088 [cs].

Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. Swe-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=VTF8yNQM66.

Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, Zhuofu Chen, Jialei Cui, Hao Ding, Mengnan Dong, Angang Du, Chenzhuang Du, Dikang Du, Yulun Du, Yu Fan, Yichen Feng, Kelin Fu, Bofei Gao, Hongcheng Gao, Peizhong Gao, Tong Gao, Xinran Gu, Longyu Guan, Haiqing Guo, Jianhang Guo, Hao Hu, Xiaoru Hao, Tianhong He, Weiran He, Wenyang He, Chao Hong, Yangyang Hu, Zhenxing Hu, Weixiao Huang, Zhiqi Huang, Zihao Huang, Tao Jiang, Zhejun Jiang, Xinyi Jin, Yongsheng Kang, Guokun Lai, Cheng Li, Fang Li, Haoyang Li, Ming Li, Wentao Li, Yanhao Li, Yiwei Li, Zhaowei Li, Zheming Li, Hongzhan Lin, Xiaohan Lin, Zongyu Lin, Chengyin Liu, Chenyu Liu, Hongzhang Liu, Jingyuan Liu, Junqi Liu, Liang Liu, Shaowei Liu, T. Y. Liu, Tianwei Liu, Weizhou Liu, Yangyang Liu, Yibo Liu, Yiping Liu, Yue Liu, Zhengying Liu, Enzhe Lu, Lijun Lu, Shengling Ma, Xinyu Ma, Yingwei Ma, Shaoguang Mao, Jie Mei, Xin Men, Yibo Miao, Siyuan Pan, Yebo Peng, Ruoyu Qin, Bowen Qu, Zeyu Shang, Lidong Shi, Shengyuan Shi, Feifan Song, Jianlin Su, Zhengyuan Su, Xinjie Sun, Flood Sung, Heyi Tang, Jiawen Tao, Qifeng Teng, Chensi Wang, Dinglu Wang, Feng Wang, Haiming Wang, Jianzhou Wang, Jiaxing Wang, Jinhong Wang, Shengjie Wang, Shuyi Wang, Yao Wang, Yejie Wang, Yiqin Wang, Yuxin Wang, Yuzhi Wang, Zhaoji Wang, Zhengtao Wang, Zhexu Wang, Chu Wei, Qianqian Wei, Wenhao Wu, Xingzhe Wu, Yuxin Wu, Chenjun Xiao, Xiaotong Xie, Weimin Xiong, Boyu Xu, Jing Xu, Jinjing Xu, L. H. Xu, Lin Xu, Suting Xu, Weixin Xu, Xinran Xu, Yangchuan Xu, Ziyao Xu, Junjie Yan, Yuzi Yan, Xiaofei Yang, Ying Yang, Zhen Yang, Zhilin Yang, Zonghan Yang, Haotian Yao, Xingcheng Yao, Wenjie Ye, Zhuorui Ye, Bohong Yin, Longhui Yu, Enming Yuan, Hongbang Yuan, Mengjie Yuan, Haobing Zhan, Dehao Zhang, Hao Zhang, Wanlu Zhang, Xiaobin Zhang, Yangkun Zhang, Yizhi Zhang, Yongting Zhang, Yu Zhang, Yutao Zhang, Yutong Zhang, Zheng Zhang, Haotian Zhao, Yikai Zhao, Huabin Zheng, Shaojie Zheng, Jianren Zhou, Xinyu Zhou, Zaida Zhou, Zhen Zhu, Weiyu Zhuang, and Xinxing Zu. Kimi K2: Open Agentic Intelligence, July 2025. URL http://arxiv.org/abs/2507.20534. arXiv:2507.20534 [cs].

Yeskendir Koishekenov, Alexandre Berard, and Vassilina Nikoulina. Memory-efficient NLLB-200: Language-specific Expert Pruning of a Massively Multilingual Machine Translation Model. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3567–3585, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.198. URL https://aclanthology.org/2023.acl-long.198/.

Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James Validad Miranda, Alisa Liu, Nouha Dziri, Xinxi Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Christopher Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. Tulu 3: Pushing frontiers in open language model post-training. In *Second Conference on Language Modeling*, 2025. URL https://openreview.net/forum?id=iluGbfHHpH.

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 *Proceedings of the Ninth International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=grwe7XHTmYb.

- Pingzhi Li, Zhenyu Zhang, Prateek Yadav, Yi-Lin Sung, Yu Cheng, Mohit Bansal, and Tianlong Chen. Merge, Then Compress: Demystify Efficient SMoE with Hints from Its Routing Policy. In *Proceedings of the Twelfth International Conference on Learning Representations*, October 2023. URL https://openreview.net/forum?id=eFWG9Cy3WK.
- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. In *The Twelfth International Conference on Learning Representations*, 2023.
- Bill Yuchen Lin, Yuntian Deng, Khyathi Chandu, Abhilasha Ravichander, Valentina Pyatkin, Nouha Dziri, Ronan Le Bras, and Yejin Choi. WildBench: Benchmarking LLMs with Challenging Tasks from Real Users in the Wild. In *Proceedings of the Thirteenth International Conference on Learning Representations*, October 2024. URL https://openreview.net/forum?id=MKEHCx25xp.
- Enshu Liu, Junyi Zhu, Zinan Lin, Xuefei Ning, Matthew B. Blaschko, Shengen Yan, Guohao Dai, Huazhong Yang, and Yu Wang. Efficient Expert Pruning for Sparse Mixture-of-Experts Language Models: Enhancing Performance and Reducing Inference Costs, July 2024a. URL http://arxiv.org/abs/2407.00945. arXiv:2407.00945 [cs].
- James Liu, Pragaash Ponnusamy, Tianle Cai, Han Guo, Yoon Kim, and Ben Athiwaratkun. Training-Free Activation Sparsity in Large Language Models, October 2024b. URL http://arxiv.org/abs/2408.14690. arXiv:2408.14690 [cs].
- Jiacheng Liu, Peng Tang, Wenfeng Wang, Yuhang Ren, Xiaofeng Hou, Pheng-Ann Heng, Minyi Guo, and Chao Li. A Survey on Inference Optimization Techniques for Mixture of Experts Models, January 2025. URL http://arxiv.org/abs/2412.14219. arXiv:2412.14219 [cs] version: 2.
- Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *Advances in Neural Information Processing Systems*, 36:21558–21572, 2023.
- Zuxin Liu, Thai Quoc Hoang, Jianguo Zhang, Ming Zhu, Tian Lan, Shirley Kokane, Juntao Tan, Weiran Yao, Zhiwei Liu, Yihao Feng, Rithesh R N, Liangwei Yang, Silvio Savarese, Juan Carlos Niebles, Huan Wang, Shelby Heinecke, and Caiming Xiong. APIGen: Automated Plpeline for generating verifiable and diverse function-calling datasets. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024c. URL https://openreview.net/forum?id=Jfg3vw2bjx.
- 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 Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 6159–6172, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10. 18653/v1/2024.acl-long.334. URL https://aclanthology.org/2024.acl-long.334/.
- Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with evol-instruct. In *Proceedings of the Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=UnUwSIgK5W.
- Meta AI Team. The Llama 4 herd: The beginning of a new era of natively multimodal AI innovation, 2025. URL https://ai.meta.com/blog/llama-4-multimodal-intelligence/.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *EMNLP*, 2018.
- ModelScope Team. EvalScope: Evaluation framework for large models, 2024. URL https://github.com/modelscope/evalscope.
- Alexandre Muzio, Alex Sun, and Churan He. SEER-MoE: Sparse Expert Efficiency through Regularization for Mixture-of-Experts, April 2024. URL http://arxiv.org/abs/2404.05089.arXiv:2404.05089 [cs].

- Neil Chowdhury, James Aung, Chan Jun Shern, Oliver Jaffe, Dane Sherburn, Giulio Starace, Evan Mays, Rachel Dias, Marwan Aljubeh, Mia Glaese, Carlos E. Jimenez, John Yang, Leyton Ho, Tejal Patwardhan, Kevin Liu, and Aleksander Madry. Introducing SWE-bench Verified, August 2024. URL https://openai.com/index/introducing-swe-bench-verified/.
- Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. Gorilla: Large language model connected with massive APIs. In *Proceedings of the Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=tBRNC6YemY.
- Shishir G Patil, Huanzhi Mao, Fanjia Yan, Charlie Cheng-Jie Ji, Vishnu Suresh, Ion Stoica, and Joseph E. Gonzalez. The berkeley function calling leaderboard (BFCL): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=2GmDdhBdDk.
- Jade Pritsker. euclaise/WritingPrompts\_curated · Datasets at Hugging Face, December 2024. URL https://huggingface.co/datasets/euclaise/WritingPrompts\_curated.
- 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. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- RedHatAI. RedHatAI/Kimi-K2-Instruct-quantized.w4a16 · Hugging Face, September 2025. URL https://huggingface.co/RedHatAI/Kimi-K2-Instruct-quantized.w4a16.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021. URL https://dl.acm.org/doi/abs/10.1145/3474381.
- Salesforce. Salesforce/xlam-function-calling-60k · Datasets at Hugging Face, May 2025. URL https://huggingface.co/datasets/Salesforce/xlam-function-calling-60k.
- Ekansh Sharma, Daniel M. Roy, and Gintare Karolina Dziugaite. The Non-Local Model Merging Problem: Permutation Symmetries and Variance Collapse, October 2024. URL http://arxiv.org/abs/2410.12766.arXiv:2410.12766 [cs].
- Noam Shazeer. GLU Variants Improve Transformer, February 2020. URL http://arxiv.org/abs/2002.05202.arXiv:2002.05202 [cs].
- Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously Large Neural Networks: The Sparsely-Gated Mixture-of-Experts Layer, January 2017. URL http://arxiv.org/abs/1701.06538.arXiv:1701.06538 [cs, stat].
- Zhi Rui Tam. theblackcat102/evol-codealpaca-v1 · Datasets at Hugging Face, July 2023. URL https://huggingface.co/datasets/theblackcat102/evol-codealpaca-v1.
- Prateek Yadav, Derek Tam, Leshem Choshen, Colin Raffel, and Mohit Bansal. TIES-merging: Resolving interference when merging models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=xtaX3WyCjl.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 Technical Report, May 2025a. URL http://arxiv.org/abs/2505.09388.arXiv:2505.09388 [cs].
- Cheng Yang, Yang Sui, Jinqi Xiao, Lingyi Huang, Yu Gong, Yuanlin Duan, Wenqi Jia, Miao Yin, Yu Cheng, and Bo Yuan. MoE-i<sup>2</sup>: Compressing mixture of experts models through inter-expert pruning and intra-expert low-rank decomposition. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 10456–10466, Miami, Florida, USA, November 2024a. Association for Computational Linguistics. doi: 10.18653/v1/2024. findings-emnlp.612. URL https://aclanthology.org/2024.findings-emnlp.612/.

 John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik R Narasimhan, and Ofir Press. SWE-agent: Agent-computer interfaces enable automated software engineering. In *Proceedings of the Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024b. URL https://openreview.net/forum?id=mXpq6ut8J3.

John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik R Narasimhan, and Ofir Press. SWE-bench/SWE-smith-trajectories · Datasets at Hugging Face, May 2025b. URL https://huggingface.co/datasets/SWE-bench/SWE-smith-trajectories.

- John Yang, Kilian Lieret, Carlos E. Jimenez, Alexander Wettig, Kabir Khandpur, Yanzhe Zhang, Binyuan Hui, Ofir Press, Ludwig Schmidt, and Diyi Yang. SWE-smith: Scaling Data for Software Engineering Agents, May 2025c. URL http://arxiv.org/abs/2504.21798. arXiv:2504.21798 [cs].
- Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. Language models are super mario: Absorbing abilities from homologous models as a free lunch. In *Forty-first International Conference on Machine Learning*, 2024.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a machine really finish your sentence? In Anna Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 4791–4800, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1472. URL https://aclanthology.org/P19-1472/.

Aohan Zeng, Xin Lv, Qinkai Zheng, Zhenyu Hou, Bin Chen, Chengxing Xie, Cunxiang Wang, Da Yin, Hao Zeng, Jiajie Zhang, Kedong Wang, Lucen Zhong, Mingdao Liu, Rui Lu, Shulin Cao, Xiaohan Zhang, Xuancheng Huang, Yao Wei, Yean Cheng, Yifan An, Yilin Niu, Yuanhao Wen, Yushi Bai, Zhengxiao Du, Zihan Wang, Zilin Zhu, Bohan Zhang, Bosi Wen, Bowen Wu, Bowen Xu, Can Huang, Casey Zhao, Changpeng Cai, Chao Yu, Chen Li, Chendi Ge, Chenghua Huang, Chenhui Zhang, Chenxi Xu, Chenzheng Zhu, Chuang Li, Congfeng Yin, Daoyan Lin, Dayong Yang, Dazhi Jiang, Ding Ai, Erle Zhu, Fei Wang, Gengzheng Pan, Guo Wang, Hailong Sun, Haitao Li, Haiyang Li, Haiyi Hu, Hanyu Zhang, Hao Peng, Hao Tai, Haoke Zhang, Haoran Wang, Haoyu Yang, He Liu, He Zhao, Hongwei Liu, Hongxi Yan, Huan Liu, Huilong Chen, Ji Li, Jiajing Zhao, Jiamin Ren, Jian Jiao, Jiani Zhao, Jianyang Yan, Jiaqi Wang, Jiayi Gui, Jiayue Zhao, Jie Liu, Jijie Li, Jing Li, Jing Lu, Jingsen Wang, Jingwei Yuan, Jingxuan Li, Jingzhao Du, Jinhua Du, Jinxin Liu, Junkai Zhi, Junli Gao, Ke Wang, Lekang Yang, Liang Xu, Lin Fan, Lindong Wu, Lintao Ding, Lu Wang, Man Zhang, Minghao Li, Minghuan Xu, Mingming Zhao, Mingshu Zhai, Pengfan Du, Qian Dong, Shangde Lei, Shangqing Tu, Shangtong Yang, Shaoyou Lu, Shijie Li, Shuang Li, Shuang-Li, Shuxun Yang, Sibo Yi, Tianshu Yu, Wei Tian, Weihan Wang, Wenbo Yu, Weng Lam Tam, Wenjie Liang, Wentao Liu, Xiao Wang, Xiaohan Jia, Xiaotao Gu, Xiaoying Ling, Xin Wang, Xing Fan, Xingru Pan, Xinyuan Zhang, Xinze Zhang, Xiuqing Fu, Xunkai Zhang, Yabo Xu, Yandong Wu, Yida Lu, Yidong Wang, Yilin Zhou, Yiming Pan, Ying Zhang, Yingli Wang, Yingru Li, Yinpei Su, Yipeng Geng, Yitong Zhu, Yongkun Yang, Yuhang Li, Yuhao Wu, Yujiang Li, Yunan Liu, Yunqing Wang, Yuntao Li, Yuxuan Zhang, Zezhen Liu, Zhen Yang, Zhengda Zhou, Zhongpei Qiao, Zhuoer Feng, Zhuorui Liu, Zichen Zhang, Zihan Wang, Zijun Yao, Zikang Wang, Ziqiang Liu, Ziwei Chai, Zixuan Li, Zuodong Zhao, Wenguang Chen, Jidong Zhai, Bin Xu, Minlie Huang, Hongning Wang, Juanzi Li, Yuxiao Dong, and Jie Tang. GLM-4.5: Agentic, Reasoning, and Coding (ARC) Foundation Models, August 2025. URL http://arxiv.org/abs/2508.06471. arXiv:2508.06471 [cs].

## A EXTENSION TO HIERARCHICAL CLUSTERING

While Theorem 1 analyses pairwise merging, practical implementations often employ hierarchical clustering to form groups of experts. Consider a cluster  $C = \{f_{i_1},...,f_{i_k}\}$  of k experts merged into a single representative  $\tilde{f}_C$ . The original contribution of this cluster can be decomposed as:

$$\sum_{j \in C} g_{i_j}(x) f_{i_j}(x) = \left(\sum_{j \in C} g_{i_j}(x)\right) \cdot \sum_{j \in C} w_j(x) f_{i_j}(x) \tag{10}$$

where  $w_j(x) = \frac{g_{i_j}(x)}{\sum_{l \in C} g_{i_l}(x)}$  are the within-cluster mixing ratios that sum to 1.

After hierarchical merging, the router must apply the *summed gate*  $\sum_{j \in C} g_{i_j}$  to a *single, static* cluster representative  $\tilde{f}_C$ , typically computed as a weighted average of the cluster members based on calibration data. This induces an irreducible error:

**Theorem 2** (Hierarchical clustering error). For a cluster C merged into  $\tilde{f}_C = \sum_{j \in C} \alpha_j f_{ij}$  with fixed weights  $\alpha_j \ge 0$ ,  $\sum_j \alpha_j = 1$ , the minimal  $L^2$  error is:

$$\min_{\{\alpha_j\}} \left\| \sum_{j \in C} g_{i_j} f_{i_j} - \left( \sum_{j \in C} g_{i_j} \right) \tilde{f}_C \right\|^2 = \mathbb{E} \left[ \left( \sum_{j \in C} g_{i_j} \right)^2 \right] \cdot \operatorname{Var}_x \left[ \sum_{j \in C} w_j(x) f_{i_j}(x) \right] \tag{11}$$

The error grows with both the cluster's total gate-value and the variance of the dynamic mixture that the cluster must approximate with a static representative.

**Implications for cluster formation.** The hierarchical error bound reveals a fundamental tension:

- Large clusters (|C| large) aggregate more gate-value  $\sum_{j \in C} g_{i_j}$ , amplifying any approximation error
- Diverse clusters (high  $\|\Delta_{ij}\|$  for  $i,j \in C$ ) increase the variance term, as the static representative must approximate a wider range of functions
- Imbalanced clustering (many singletons, few mega-clusters) combines the worst aspects: mega-clusters suffer severe collapse while singletons provide minimal compression

Distance metrics like Euclidean distance that consider magnitude can exacerbate these issues by creating clusters based on norm similarity rather than functional role, potentially grouping experts with different specializations but similar scales. The resulting mega-clusters force the router to apply a single control signal to what were previously dozens of independently modulated experts, explaining the catastrophic functional collapse observed empirically in late layers where  $Var[w_i(x)]$  is highest.

## B ADDITIONAL EMPIRICAL EVIDENCE FOR LOSS OF INDEPENDENT CONTROL

In Figure 1a, Qwen3's layer 0 exemplifies the contraction of the functional output space by merging in early layers. The original 128 experts span from -0.4 to 1.0 along PC1, pruning maintains this full range with 64 experts, while merging contracts the distribution to approximately [-0.2,0.3], a 5-fold reduction. This contraction is dramatic in late layers, where experts are more specialized. As depicted in Figure A4f, the original 15 experts of Llama-4's layer 47 occupy a vast, multi-modal space spanning PC1 coordinates from -800 to 600. Pruning preserves this remarkable diversity, with the 8 surviving experts distributed across the same multi-modal regions. However, merging induces a catastrophic collapse to a tiny cluster around coordinates (200,0), representing nearly two orders of magnitude reduction in functional diversity. This pattern intensifies with the number of experts: Qwen3's layer 47 (Figure 1b) shows the most severe collapse, with 128 original experts spanning PC1 from -200 to 300 reduced to a minute region after merging, while its 64 pruned experts maintain the original distribution's full breadth.

**Manifold geometry preservation** Across all models and layers, we observe a fundamental geometric principle: pruning preserves the topology of the functional manifold while merging fundamentally alters it. This distinction is most clearly visible in ERNIE's representations (Figures A4a and A4b). In layer 1, the original 64 routed experts plus 2 shared experts form a characteristic curved structure with several

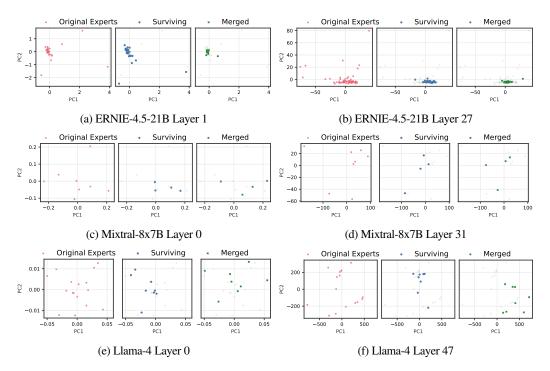


Figure A4: (a,c,e) Functional subspace (PCA) for early SMoE layers. Pruning (blue) preserves the manifold geometry; merging (green) collapses it toward the centre. (b,d,f) Functional subspace (PCA) for late MoE layers.

outliers representing specialized experts. After pruning, the red points precisely overlay the gray ghost of the original distribution, including the outlier positions, demonstrating that each surviving expert maintains its exact functional role. The merged configuration, however, shows all experts collapsed into a tight cluster at the distribution's centre, eliminating both the outliers and the manifold's curvature.

The preservation of manifold geometry under pruning reflects the mathematical structure of the operation: the pruned hypothesis class is a coordinate subspace of the original, with the router maintaining independent control over each surviving expert. The geometric collapse under merging visualizes the loss of independent control when gates  $g_i$  and  $g_j$  are tied into their sum  $(g_i+g_j)$ , the router can no longer independently modulate the two underlying functions, forcing the model to approximate the dynamic mixture  $r(x)f_i(x)+(1-r(x))f_j(x)$  with a static expert  $\tilde{f}_{\alpha}$ .

Mixtral, with only 8 experts, provides an interesting edge case (Figures A4c and A4d). Even with fewer experts, the same geometric principles apply. Pruning maintains the convex hull of the original distribution while merging contracts it. The less dramatic collapse compared to models with more experts suggests that with fewer experts, each must remain more general, leading to lower  $\|\Delta_{ij}\|^2$  (expert gap) and lower Var[r(x)] (policy variability), both factors in our irreducible error bound.

## C EVALUATION DETAILS

Multiple choice (MC) evaluation. Following Chen et al. (2025), our MC benchmarks include: AI2 Reasoning Challenge (ARC-c & ARC-e) (Clark et al., 2018), BoolQ (Clark et al., 2019), HellaSwag (Zellers et al., 2019), MMLU (Hendrycks et al., 2021a), OpenBookQA (OBQA) (Mihaylov et al., 2018), Recognizing Textual Entailment Challenge (RTE) (Bentivogli et al., 2009), and WinoGrande (WinoG.) (Sakaguchi et al., 2021). We evaluate the models in the zero-shot setting using the standard log-likelihood approach with Im-eval-harness (Gao et al., 2023). We report byte-length normalized accuracies for ARC-c, ARC-e, HellaSwag, and OBQA<sup>1</sup>.

<sup>&</sup>lt;sup>1</sup>Reported as the acc\_norm field in the EleutherAI evaluation harness outputs. See Gao (2021) for details.

 Coding evaluation. For code generation, all models are evaluated on EvalPlus (Liu et al., 2023) and 182 LiveCodeBench (Jain et al., 2025) questions collected between January and April 2025. We extend the original source code for these benchmarks to evaluate our models. We additionally evaluate Kimi-K2-Instruct-W4A16 and Qwen3-Coder-480B on the agentic coding benchmark SWE-Bench (Jimenez et al., 2024) and tool-calling benchmark BFCLv3 (Patil et al., 2025). For BFCLv3, we use the original Gorilla framework for evaluating our models (Patil et al., 2024).

For SWE-Bench evaluation, we run our compressed models with the mini-SWE-agent scaffolding (Yang et al., 2024b) and report the score on the SWE-Bench Verified test set (Neil Chowdhury et al., 2024). We use 4,096 and 16,384 as the maximum number of output tokens for evaluating Qwen3-Coder-480B and Kimi-K2-Instruct-W4A16 on SWE-Bench, respectively. The input context length for both models is limited to 65,536. We do not limit the number of turns in mini-SWE-agent flow, but restart the rollout in cases where the model could not generate a valid patch (that is, in the case when the output of the final turn does not contain a diff-git substring). We set the maximum number of restarts to 20, which we found to be sufficient to generate patches for all samples with pruned models, unless the model produces degenerate responses like repeating strings. We use the cloud-based evaluation provided with the sb-cli tool to get the final scores for all evaluated models.

Math and creative writing evaluation. Mathematical reasoning is assessed on GSM8K (Cobbe et al., 2021) and MATH-500 (Hendrycks et al., 2021b; Lightman et al., 2023) benchmarks using the evalscope (ModelScope Team, 2024) framework. To assess creative writing, we use 146 creative writing prompts sampled from WildBench (Lin et al., 2024) with GPT-40 used as the judge to evaluate the model responses. We report normalized scores using the WildBench rubric.

**Generation configuration.** For models with  $\leq 110B$  parameters, we use greedy sampling (i.e, temperature = 0.0) to evaluate code generation and math reasoning. For creative writing we use the default temperature, top-P, and top-K settings for each respective model. The maximum number of output tokens is extended to 16,384 for all generative tasks to account for the verbosity of some models. For hybrid reasoning models such as Qwen3-30B-A3B, we disable reasoning on all tasks by setting enable\_thinking=False in the chat template.

For larger models with  $\geq 110B$  parameters, we use greedy sampling for EvalPlus, SWE-Bench, and BFCLv3. On LiveCodeBench, Qwen3-Coder-480B and Kimi-K2 are evaluated with default sampling parameters and greedy sampling, respectively. We report the mean and standard deviation for Qwen3-Coder-480B on LiveCodeBench over five random seeds. We use a repetition penalty of 1.05 for all large model evaluations. For EvalPlus we use 768 as the maximum number of output tokens and 16,384 for LiveCodeBench. For BFCLv3 we set the maximum number of output tokens to 4,096.

The Kimi-K2-Instruct-W4A16 model used throughout this study is an INT4 weight-quantized version of Kimi-K2-Instruct released by RedHatAI (2025).

## D BASELINE METHODS

The following formally describes the baselines compression methods we consider.

**Notation.** Let  $\mathcal{X}_{cal}$  be a calibration dataset. Consider a SMoE model with n layers,  $L_n$ , K experts per layer  $f_1, \ldots, f_K$ , each a function  $f_k : \mathbb{R}^d \to \mathbb{R}^d$ , and a router producing non-negative gates  $\mathbf{g}(x) = (g_1(x), \ldots, g_K(x)) \in \Delta^{K-1}$ . The output of layer  $L_n$  is

$$h_n = \sum_{i=1}^{K} g_i(x) f_i(x).$$

The expert usage frequency,  $\nu_i$ , for expert  $f_i$  is the number of tokens in  $\mathcal{X}_{cal}$  for which  $f_i$  is activated

$$\nu_i = |\mathcal{X}_i|,$$

where  $\mathcal{X}_i = \{x \in \mathcal{X}_{cal} \mid i \in \text{TopK}(\mathbf{g}(x))\}.$ 

Given saliency scores,  $\mathbf{S} \in \mathbb{R}^K$ , pruning removes experts with the minimum saliency score. For merging, we first cluster experts based on their pairwise distances,  $\mathbf{D} \in \mathbb{R}^{K \times K}$ , and then merge the parameters of experts contained within each cluster.

**Frequency-based pruning.** The frequency-based pruning saliency criterion prunes experts with the lowest usage frequency across the calibration dataset. The saliency of  $f_i$  is simply  $S_i = \nu_i$ .

**EAN pruning.** EAN pruning introduced by Jaiswal et al. (2025) accumulates the activation norm of each expert across tokens for which the expert is activated. The saliency of  $f_i$  is

$$S_i = \sum_{x \in \mathcal{X}_i} \|f_i(x)\|_2. \tag{12}$$

**M-SMoE merging.** Proposed by Li et al. (2023), M-SMoE first uses weight-matching (Ainsworth et al., 2023) to find a permutation matrix  $\mathbf{P_j}$  which aligns expert  $f_j$  to expert  $f_i$ . In the models we study, each expert is a two-layer feed-forward SwiGLU block (Shazeer, 2020) with up, gate, and down projections:  $f_j = \{W_{up}^{(j)}, W_{gate}^{(j)}, W_{down}^{(j)}\}$ . The permutation matrix is applied to the intermediate dimension of the experts such that the expert outputs are invariant to the transformation

$$W_{up}^{\prime(j)}\!=\!W_{up}^{(j)}\mathbf{P}_j, \hspace{1cm} W_{gate}^{\prime(j)}\!=\!W_{gate}^{(j)}\mathbf{P}_j, \hspace{1cm} W_{down}^{\prime(j)}\!=\!\mathbf{P}_j^TW_{down}^{(j)}.$$

The permuted expert is defined as  $\tilde{f}_j = \{W_{up}^{\prime(j)}, W_{gate}^{\prime(j)}, W_{down}^{\prime(j)}\}$ .

To initialize the expert clusters, M-SMoE identifies the set of m dominant experts  $\mathbb{F}_{dom}$ , as the experts across all layers with the highest usage frequency  $\nu$ . The pairwise expert distance is based on the cosine distance of the router gate-values measured on the calibration dataset

$$D_{i,j} = \frac{1}{|\mathcal{X}_{cal}|} \sum_{x \in \mathcal{X}_{cal}} 1 - \frac{g_i(x) \cdot g_j(x)}{\|g_i(x)\| \|g_j(x)\|}.$$
 (13)

Non-dominant expert j is clustered by selecting the dominant expert with the smallest pairwise distance

$$i^* = \underset{i \in \mathbb{F}_{dom}}{\operatorname{argmin}} D_{i,j}$$

The merged expert  $f_{\alpha}$  is created by calculating the frequency-weighted average of the permuted parameters, W', of all experts in the cluster  $\mathbb{C}_{\alpha}$ 

$$\tilde{W}_a = \frac{\sum_{i \in \mathbb{C}_\alpha} \nu_i W_i'}{\sum_{i \in \mathbb{C}_\alpha \nu_i}}.$$
(14)

**HC-SMoE merging.** Chen et al. (2025) clusters experts based on their *representative vectors*,  $A_i$ , defined as the average activation across every token in the calibration dataset

$$A_i := \mathbb{E}_{x \sim \mathcal{X}_{cal}}[f_i(x)] = \frac{1}{|\mathcal{X}_{cal}|} \sum_{x \in \mathcal{X}_{cal}} f_i(x).$$

The expert pairwise distance is defined as the cosine distance between representative vectors

$$D_{i,j} = 1 - \frac{A_i \cdot A_j}{\|A_i\| \|A_j\|}.$$
 (15)

Clusters are formed using hierarchical agglomerative clustering with average linkage criterion. We start by initializing each expert as a singleton cluster. At every iteration, the closest pair of clusters,  $\mathbb{C}_i^*, \mathbb{C}_j^*$  are joined and the pairwise distances updated as the average of the constituents

$$i^*, j^* = \underset{i,j}{\operatorname{argmin}} D_{i,j}, \qquad \qquad \mathbb{C}_{\alpha} = \mathbb{C}_{i^*} \cup \mathbb{C}_{j^*}, \qquad \qquad D_{a,k} = \frac{\sum_{i \in \mathbb{C}_{\alpha}} D_{i,k}}{|\mathbb{C}_{\alpha}|}.$$

The clusters are merged with equation 14.

### E ADDITIONAL RESULTS

Table A4 shows the full suite of MC question answering benchmarks and the average result across all models and methods. Table A5 tabulates code generation accuracy of compressed SMoE models calibrated

Table A4: Detailed benchmark results for multiple-choice QA tasks.

Model	Compression	Technique	Method	ARC-c	ARC-e	BoolQ	Hellaswag	MMLU	OBQA	RTE	WinoG.	MC Avg
	Baseline			0.564	0.782	0.873	0.813	0.737	0.462	0.812	0.724	0.721
		Merging	M-SMoE HC-SMoE	$0.434 \pm 0.006$ $0.506 \pm 0.000$	$0.652 \pm 0.008$ $0.717 \pm 0.001$	$0.846 \pm 0.001$ $0.849 \pm 0.001$	$0.597 \pm 0.002$ $0.714 \pm 0.001$	$0.591 \pm 0.001$ $0.652 \pm 0.002$	$0.350 \pm 0.006$ $0.371 \pm 0.002$	$0.819 \pm 0.010$ $0.799 \pm 0.002$	$0.655 \pm 0.003$ $0.674 \pm 0.004$	$0.618 \pm 0.00$ $0.660 \pm 0.00$
ERNIE-4.5-21B- A3B-PT	25%	Pruning	Frequency EAN REAP	$0.486 \pm 0.004$ $0.498 \pm 0.005$ $0.527 \pm 0.004$	$0.711 \pm 0.000$ $0.713 \pm 0.002$ $0.759 \pm 0.002$	$0.852 \pm 0.004$ $0.863 \pm 0.002$ $0.857 \pm 0.003$	$0.675 \pm 0.003$ $0.717 \pm 0.004$ $0.717 \pm 0.003$	$0.628 \pm 0.003$ $0.625 \pm 0.001$ $0.644 \pm 0.001$	$0.373 \pm 0.003$ $0.405 \pm 0.011$ $0.409 \pm 0.009$	$0.780 \pm 0.006$ $0.811 \pm 0.009$ $0.756 \pm 0.008$	$0.676 \pm 0.005$ $0.702 \pm 0.005$ $0.690 \pm 0.001$	$0.648 \pm 0.00$ $0.667 \pm 0.00$ $0.670 \pm 0.00$
		Merging	M-SMoE HC-SMoE	0.294 ± 0.033 0.411 ± 0.003	0.452 ± 0.040 0.641 ± 0.002	0.764 ± 0.010 0.822 ± 0.001	0.341 ± 0.011 0.523 ± 0.001	$0.385 \pm 0.001$ $0.495 \pm 0.002$	$0.270 \pm 0.004$ $0.330 \pm 0.005$	0.687 ± 0.017 0.742 ± 0.011	0.529 ± 0.010 0.587 ± 0.009	$0.465 \pm 0.0$ $0.569 \pm 0.0$
	50%	Pruning	Frequency EAN REAP	$0.400 \pm 0.002$ $0.417 \pm 0.005$ $0.417 \pm 0.009$	$0.584 \pm 0.006$ $0.633 \pm 0.005$ $0.626 \pm 0.007$	$0.830 \pm 0.001$ $0.830 \pm 0.003$ $0.803 \pm 0.006$	$0.522 \pm 0.003$ $0.572 \pm 0.001$ $0.556 \pm 0.003$	$0.506 \pm 0.006$ $0.509 \pm 0.002$ $0.505 \pm 0.003$	$0.303 \pm 0.004$ $0.336 \pm 0.003$ $0.325 \pm 0.006$	$0.758 \pm 0.004$ $0.785 \pm 0.014$ $0.775 \pm 0.014$	$0.625 \pm 0.004$ $0.626 \pm 0.003$ $0.623 \pm 0.008$	0.566 ± 0.0 0.589 ± 0.0 0.579 ± 0.0
	Baseline			0.563	0.790	0.887	0.778	0.779	0.454	0.816	0.702	0.721
		Merging	M-SMoE HC-SMoE	$0.357 \pm 0.006$ $0.478 \pm 0.006$	$0.519 \pm 0.003$ $0.722 \pm 0.006$	$0.843 \pm 0.006 \\ 0.863 \pm 0.003$	$0.529 \pm 0.002 \\ 0.714 \pm 0.000$	$0.536 \pm 0.004 \\ 0.684 \pm 0.002$	$\begin{array}{c} 0.310 \pm 0.005 \\ 0.417 \pm 0.001 \end{array}$	$0.735 \pm 0.027$ $0.805 \pm 0.004$	$0.635 \pm 0.005 \\ 0.710 \pm 0.004$	0.558 ± 0.0 0.674 ± 0.0
Qwen3-30B-A3B	25%	Pruning	Frequency EAN REAP	$0.401 \pm 0.011$ $0.406 \pm 0.007$ $0.481 \pm 0.005$	$0.600 \pm 0.016$ $0.603 \pm 0.014$ $0.720 \pm 0.005$	$\begin{array}{c} 0.847 \pm 0.003 \\ 0.847 \pm 0.005 \\ 0.852 \pm 0.003 \end{array}$	$0.593 \pm 0.005$ $0.607 \pm 0.006$ $0.706 \pm 0.006$	$0.600 \pm 0.004$ $0.600 \pm 0.002$ $0.674 \pm 0.002$	$\begin{array}{c} 0.342 \pm 0.012 \\ 0.337 \pm 0.003 \\ 0.405 \pm 0.005 \end{array}$	$\begin{array}{c} 0.781 \pm 0.002 \\ 0.764 \pm 0.002 \\ 0.813 \pm 0.006 \end{array}$	$0.637 \pm 0.005$ $0.660 \pm 0.009$ $0.701 \pm 0.008$	0.600 ± 0.0 0.603 ± 0.0 0.669 ± 0.0
	50%	Merging	M-SMoE HC-SMoE	$\begin{array}{c} 0.278 \pm 0.003 \\ 0.368 \pm 0.002 \end{array}$	$\begin{array}{c} 0.402 \pm 0.003 \\ 0.593 \pm 0.003 \end{array}$	$\begin{array}{c} 0.753 \pm 0.004 \\ 0.740 \pm 0.003 \end{array}$	$\begin{array}{c} 0.399 \pm 0.002 \\ 0.473 \pm 0.002 \end{array}$	$\begin{array}{c} 0.366 \pm 0.004 \\ 0.516 \pm 0.003 \end{array}$	$\begin{array}{c} 0.278 \pm 0.002 \\ 0.301 \pm 0.007 \end{array}$	$\begin{array}{c} 0.586 \pm 0.014 \\ 0.724 \pm 0.004 \end{array}$	$\begin{array}{c} 0.546 \pm 0.004 \\ 0.620 \pm 0.005 \end{array}$	0.451 ± 0.0 0.542 ± 0.0
	30%	Pruning	Frequency EAN REAP	$0.285 \pm 0.001$ $0.296 \pm 0.006$ $0.344 \pm 0.004$	$0.424 \pm 0.002$ $0.426 \pm 0.009$ $0.504 \pm 0.008$	$0.779 \pm 0.003$ $0.759 \pm 0.007$ $0.745 \pm 0.005$	$0.458 \pm 0.003$ $0.471 \pm 0.002$ $0.489 \pm 0.013$	$\begin{array}{c} 0.397 \pm 0.002 \\ 0.443 \pm 0.001 \\ 0.507 \pm 0.005 \end{array}$	$0.286 \pm 0.004$ $0.291 \pm 0.009$ $0.311 \pm 0.003$	$\begin{array}{c} 0.659 \pm 0.012 \\ 0.668 \pm 0.020 \\ 0.625 \pm 0.031 \end{array}$	14 0.710 ± 0.004 12 0.6637 ± 0.005 12 0.660 ± 0.009 16 0.701 ± 0.008 14 0.546 ± 0.004 14 0.546 ± 0.004 14 0.546 ± 0.004 15 0.589 ± 0.009 16 0.701 ± 0.008 17 0.655 ± 0.004 17 0.655 ± 0.004 17 0.655 ± 0.004 18 0.732 ± 0.007 19 0.655 ± 0.004 10 0.589 ± 0.009 10 0.589 ± 0.009 10 0.589 ± 0.009 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005 10 0.589 ± 0.005	$0.483 \pm 0.0$ $0.493 \pm 0.0$ $0.518 \pm 0.0$
	Baseline			0.650	0.842	0.887	0.861	0.691	0.496	0.722	0.740	0.736
Mixtral-8x7B- Instruct-v0.1	25%	Merging	M-SMoE HC-SMoE	$\begin{array}{c} 0.532 \pm 0.004 \\ 0.590 \pm 0.004 \end{array}$	$\begin{array}{c} 0.769 \pm 0.007 \\ 0.797 \pm 0.004 \end{array}$	$\begin{array}{c} 0.847 \pm 0.001 \\ 0.869 \pm 0.003 \end{array}$	$\begin{array}{c} 0.747 \pm 0.002 \\ 0.835 \pm 0.002 \end{array}$	$\begin{array}{c} 0.553 \pm 0.001 \\ 0.626 \pm 0.000 \end{array}$	$\begin{array}{c} 0.429 \pm 0.008 \\ 0.482 \pm 0.004 \end{array}$	$\begin{array}{c} 0.632 \pm 0.010 \\ 0.703 \pm 0.012 \end{array}$	$0.731 \pm 0.007$	$0.646 \pm 0.0$ $0.704 \pm 0.0$
	2370	Pruning	Frequency EAN REAP	$\begin{array}{c} 0.616 \pm 0.014 \\ 0.607 \pm 0.004 \\ 0.611 \pm 0.003 \end{array}$	$\begin{array}{c} 0.826 \pm 0.007 \\ 0.831 \pm 0.001 \\ 0.825 \pm 0.001 \end{array}$	$\begin{array}{c} 0.875 \pm 0.001 \\ 0.884 \pm 0.001 \\ 0.874 \pm 0.002 \end{array}$	$\begin{array}{c} 0.825 \pm 0.002 \\ 0.836 \pm 0.001 \\ 0.830 \pm 0.002 \end{array}$	$\begin{array}{c} 0.637 \pm 0.003 \\ 0.646 \pm 0.002 \\ 0.643 \pm 0.001 \end{array}$	$\begin{array}{c} 0.451 \pm 0.003 \\ 0.484 \pm 0.005 \\ 0.475 \pm 0.006 \end{array}$	$\begin{array}{c} 0.706 \pm 0.017 \\ 0.700 \pm 0.004 \\ 0.761 \pm 0.002 \end{array}$	$\begin{array}{c} 0.732 \pm 0.004 \\ 0.718 \pm 0.001 \end{array}$	$0.704 \pm 0.0$ $0.715 \pm 0.0$ $0.717 \pm 0.0$
	50%	Merging	M-SMoE HC-SMoE	$\begin{array}{c} 0.446 \pm 0.005 \\ 0.539 \pm 0.003 \end{array}$	$\begin{array}{c} 0.700 \pm 0.001 \\ 0.759 \pm 0.000 \end{array}$	$\begin{array}{c} 0.788 \pm 0.003 \\ 0.851 \pm 0.001 \end{array}$	$\begin{array}{c} 0.630 \pm 0.002 \\ 0.791 \pm 0.001 \end{array}$	$\begin{array}{c} 0.430 \pm 0.001 \\ 0.543 \pm 0.000 \end{array}$	$\begin{array}{c} 0.386 \pm 0.003 \\ 0.442 \pm 0.000 \end{array}$	$\begin{array}{c} 0.570 \pm 0.000 \\ 0.700 \pm 0.004 \end{array}$	$0.712 \pm 0.002$	$0.568 \pm 0.0$ $0.667 \pm 0.0$
	30 10	Pruning	Frequency EAN REAP	$0.541 \pm 0.004$ $0.551 \pm 0.014$ $0.544 \pm 0.005$	$0.781 \pm 0.003$ $0.774 \pm 0.008$ $0.785 \pm 0.005$	$0.824 \pm 0.013$ $0.859 \pm 0.004$ $0.837 \pm 0.003$	$\begin{array}{c} 0.759 \pm 0.002 \\ 0.794 \pm 0.002 \\ 0.778 \pm 0.002 \end{array}$	$0.516 \pm 0.002$ $0.550 \pm 0.006$ $0.554 \pm 0.001$	$\begin{array}{c} 0.411 \pm 0.006 \\ 0.452 \pm 0.014 \\ 0.462 \pm 0.005 \end{array}$	$0.708 \pm 0.023$ $0.717 \pm 0.023$ $0.715 \pm 0.013$	4 0.712 ± 0.002 3 0.650 ± 0.005 3 0.693 ± 0.008 3 0.679 ± 0.005 0.692	$0.649 \pm 0.0$ $0.674 \pm 0.0$ $0.669 \pm 0.0$
	Baseline			0.627	0.848	0.879	0.823	0.803	0.462	0.765		0.738
	25%	Merging	M-SMoE HC-SMoE	0.573 0.588	0.802 0.814	0.872 0.876	0.752 0.779	0.719 0.720	0.434 0.424	0.769 0.729	0.695	0.699 0.703
Llama-4-Scout- 17B-16E- Instruct	23 10	Pruning	Frequency EAN REAP	0.584 0.582 0.594	0.817 0.816 0.830	0.876 0.872 0.872	0.779 0.777 0.788	0.733 0.735 0.756	0.438 0.446 0.452	0.773 0.791 0.769	0.679	0.711 0.712 <b>0.718</b>
	50%	Merging	M-SMoE HC-SMoE	0.498 0.526	0.717 0.781	0.856 0.862	0.676 0.718	0.609 0.628	0.388 0.386	0.787 0.726		0.649 0.661
	30%	Pruning	Frequency EAN REAP	0.518 0.510 0.561	0.734 0.750 0.802	0.860 0.857 0.869	0.704 0.712 0.745	0.652 0.650 0.682	0.398 0.398 0.432	0.765 0.762 0.762	0.662	0.661 0.663 0.689
	Baseline			0.619	0.825	0.882	0.858	0.789	0.478	0.747	0.776	0.747
		Merging	M-SMoE HC-SMoE	0.429 0.577	0.651 0.782	0.808 0.860	0.671 0.815	0.578 0.722	0.362 0.458	0.578 0.668	4 0.623 ± 0.008 0.702 7 0.635 ± 0.005 4 0.710 ± 0.004 4 0.710 ± 0.004 4 0.710 ± 0.004 4 0.546 ± 0.004 4 0.546 ± 0.004 4 0.546 ± 0.004 4 0.546 ± 0.004 4 0.546 ± 0.004 4 0.546 ± 0.004 4 0.546 ± 0.004 4 0.546 ± 0.004 4 0.546 ± 0.004 4 0.546 ± 0.004 0 0.656 ± 0.004 0 0.656 ± 0.004 0 0.732 ± 0.005 4 0.732 ± 0.005 4 0.732 ± 0.005 0 0.665 ± 0.004 0 0.666 ± 0.006 0 0.671 0 0.692 0 0.671 0 0.692 0 0.671 0 0.692 0 0.671 0 0.693 0 0.694 0 0.755 0 0.695 0 0.691 0 0.697 0 0.683 0 0.666 0 0.667 0 0.667 0 0.669 0 0.671 0 0.691 0 0.776 0 0.692 0 0.671 0 0.691 0 0.679 0 0.683 0 0.6661 0 0.679 0 0.683 0 0.6662 0 0.6601 0 0.776 0 0.695 0 0.755 0 0.725 0 0.730 0 0.724 0 0.631 0 0.632 0 0.640 0 0.771 0 0.632 0 0.6640 0 0.7701 0 0.631 0 0.635 0 0.635 0 0.635 0 0.635 0 0.635 0 0.635 0 0.635 0 0.635 0 0.635 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631 0 0.631	0.596 <b>0.704</b>
GLM-4.5-Air	25%	Pruning	Frequency EAN REAP	0.493 0.492 0.555	0.715 0.705 0.756	0.827 0.805 0.813	0.732 0.736 0.796	0.653 0.656 0.701	0.422 0.368 0.434	0.614 0.603 0.643	0.730	0.648 0.637 <u>0.678</u>
		Merging	M-SMoE HC-SMoE	0.291 0.428	0.452 0.671	0.693 0.761	0.433 0.590	0.382 0.524	0.266 0.318	0.484 0.603	0.551	0.444 <b>0.564</b>
	50%	Pruning	Frequency EAN REAP	0.334 0.358 0.427	0.535 0.530 0.604	0.767 0.682 0.662	0.566 0.573 0.642	0.478 0.489 0.569	0.288 0.300 0.318	0.567 0.516 0.606	0.635	0.521 0.511 0.559
	Baseline			0.644	0.822	0.906	0.841	0.850	0.468	0.751		0.750
Qwen3-Coder- 480B-A35B-	25%	Pruning	Frequency EAN	0.443 0.555	0.673 0.766	0.845 0.891	0.651 0.769	0.621 0.795	0.280 0.404	0.704 0.747	0.691	0.606 0.702
Instruct-FP8	50%	Pruning	Frequency EAN	0.635 0.314 0.402	0.824 0.470 0.596	0.900 0.791 0.858	0.841 0.502 0.629	0.836 0.451 0.615	0.466 0.262 0.216	0.754 0.679 0.744	0.580 0.666	0.748 0.506 0.591
			REAP	0.546	0.772	0.872	0.756	0.696	0.430	0.762		0.692
	Baseline		P	0.712	0.879	0.913	0.765	0.872	0.504	0.783		0.780
Kimi-K2- Instruct- W4A16	25%	Pruning	Frequency EAN REAP	0.518 0.615 0.671	0.771 0.819 0.854	0.825 0.893 0.907	0.787 0.843 0.860	0.242 0.500 0.809	0.420 0.446 0.470	0.653 0.762 0.805		0.604 0.703 <b>0.773</b>
W4A16	50%	Pruning	Frequency EAN	0.285 0.426	0.498	0.620 0.863	0.436 0.663	0.241	0.314	0.617 0.726	0.500	0.439 0.587

on evol-codealpaca. Eval+ is the average of MBPP, MBPP+, HumanEval (HE), HE+. The *Code Avg* column is the average of Eval+ and LiveCodeBench (LiveCode). Table A6 summarizes the accuracy of the various compression methods studied when calibrated with the c4 dataset on coding and MC benchmarks. Notably, while the MC performance is generally slightly higher than models calibrated on evol-codealpaca, the resulting code generation quality is abysmal, with most models failing to generate coherent output.

Figure A5 plots non-agentic coding and MC accuracy versus compressed model size. Figure A6a depict the proportion of singleton clusters for HC-SMoE and M-SMoE. Figure A6b plots accuracy vs. maximum cluster sizes when the maximum cardinality of clusters is restricted. Figures A7 and A8 show the importance of using domain-specific calibration data, particularly at high compression ratios.

Table A5: Detailed benchmark results for non-agentic code generation tasks. Eval+ is the average of MBPP, MBPP+, HE, HE+. The Code Avg column is the average of Eval+ and LiveCodeBench (LiveCode).

Model	Compression	Technique	Method	HE	HE+	MBPP	MBPP+	Eval+	LiveCode	Code Avg
	Baseline	• •		0.902	0.866	0.910	0.765	0.861	0.231	0.546
		Merging	M-SMoE HC-SMoE	$0.774 \pm 0.011$ $0.837 \pm 0.007$	$0.730 \pm 0.009$ $0.805 \pm 0.000$	$0.768 \pm 0.015$ $0.827 \pm 0.003$	$0.647 \pm 0.017$ $0.696 \pm 0.008$	$0.730 \pm 0.005$ $0.791 \pm 0.004$	$0.194 \pm 0.022$ $0.207 \pm 0.008$	$0.462 \pm 0.01$ $0.499 \pm 0.002$
ERNIE-4.5-21B- A3B-PT	25%	Pruning	Frequency EAN REAP	$0.890 \pm 0.006$ $0.890 \pm 0.006$ $0.892 \pm 0.009$	$0.846 \pm 0.009$ $0.848 \pm 0.011$ $0.854 \pm 0.012$	$0.837 \pm 0.010$ $0.840 \pm 0.006$ $0.876 \pm 0.000$	$0.709 \pm 0.010$ $0.727 \pm 0.004$ $0.738 \pm 0.003$	$0.820 \pm 0.006$ $0.826 \pm 0.004$ $0.840 \pm 0.005$	$0.151 \pm 0.096$ $0.161 \pm 0.111$ $0.167 \pm 0.124$	$0.486 \pm 0.04$ $0.494 \pm 0.05$ $0.504 \pm 0.06$
		Merging	M-SMoE HC-SMoE	$0.104 \pm 0.022$ $0.425 \pm 0.004$	$0.100 \pm 0.029$ $0.404 \pm 0.007$	$0.239 \pm 0.036$ $0.608 \pm 0.018$	$0.207 \pm 0.040$ $0.511 \pm 0.011$	$0.162 \pm 0.012$ $0.487 \pm 0.008$	$0.024 \pm 0.008$ $0.082 \pm 0.015$	$0.093 \pm 0.003$ $0.285 \pm 0.003$
	50%	Pruning	Frequency EAN REAP	$0.699 \pm 0.031$ $0.675 \pm 0.019$ $0.797 \pm 0.009$	$0.640 \pm 0.022$ $0.642 \pm 0.009$ $0.764 \pm 0.007$	$0.696 \pm 0.014$ $0.713 \pm 0.015$ $0.767 \pm 0.017$	$0.584 \pm 0.006$ $0.591 \pm 0.016$ $0.644 \pm 0.013$	$0.655 \pm 0.015$ $0.655 \pm 0.014$ $0.743 \pm 0.008$	$0.083 \pm 0.066$ $0.112 \pm 0.064$ $0.137 \pm 0.119$	$0.369 \pm 0.022$ $0.384 \pm 0.032$ $0.440 \pm 0.064$
	Baseline			0.927	0.884	0.881	0.743	0.859	0.302	0.581
		Merging	M-SMoE HC-SMoE	$0.878 \pm 0.012$ $0.866 \pm 0.011$	$\begin{array}{c} 0.833 \pm 0.007 \\ 0.805 \pm 0.016 \end{array}$	$0.849 \pm 0.007$ $0.832 \pm 0.006$	$0.728 \pm 0.007$ $0.698 \pm 0.005$	$\begin{array}{c} 0.822 \pm 0.004 \\ 0.800 \pm 0.004 \end{array}$	$\begin{array}{c} 0.293 \pm 0.017 \\ 0.258 \pm 0.000 \end{array}$	$0.558 \pm 0.00$ $0.529 \pm 0.00$
Qwen3-30B-A3B	25%	Pruning	Frequency EAN REAP	$ \begin{vmatrix} 0.921 \pm 0.006 \\ 0.909 \pm 0.006 \\ 0.917 \pm 0.007 \end{vmatrix} $	$\begin{array}{c} 0.874 \pm 0.007 \\ 0.864 \pm 0.004 \\ 0.876 \pm 0.004 \end{array}$	$\begin{array}{c} 0.868 \pm 0.000 \\ 0.859 \pm 0.009 \\ 0.853 \pm 0.002 \end{array}$	$\begin{array}{c} 0.735 \pm 0.003 \\ 0.729 \pm 0.008 \\ 0.727 \pm 0.006 \end{array}$	$\begin{array}{c} 0.849 \pm 0.004 \\ 0.840 \pm 0.004 \\ 0.843 \pm 0.002 \end{array}$	$\begin{array}{c} 0.302 \pm 0.011 \\ 0.311 \pm 0.018 \\ 0.308 \pm 0.015 \end{array}$	$egin{array}{l} \textbf{0.576} \pm \textbf{0.00} \\ \textbf{0.576} \pm \textbf{0.01} \\ 0.575 \pm 0.00 \end{array}$
	50%	Merging	M-SMoE HC-SMoE	$\begin{array}{c} 0.687 \pm 0.013 \\ 0.577 \pm 0.023 \end{array}$	$\begin{array}{c} 0.638 \pm 0.004 \\ 0.541 \pm 0.013 \end{array}$	$\begin{array}{c} 0.618 \pm 0.004 \\ 0.631 \pm 0.010 \end{array}$	$\begin{array}{c} 0.541 \pm 0.007 \\ 0.546 \pm 0.004 \end{array}$	$\begin{array}{c} 0.621 \pm 0.006 \\ 0.574 \pm 0.010 \end{array}$	$\begin{array}{c} 0.205 \pm 0.019 \\ 0.185 \pm 0.018 \end{array}$	$0.413 \pm 0.00$ $0.379 \pm 0.00$
	30%	Pruning	Frequency EAN REAP	$0.787 \pm 0.016$ $0.886 \pm 0.025$ $0.919 \pm 0.007$	$0.756 \pm 0.022$ $0.837 \pm 0.020$ $0.870 \pm 0.004$	$0.692 \pm 0.016$ $0.798 \pm 0.006$ $0.805 \pm 0.009$	$0.579 \pm 0.016$ $0.669 \pm 0.008$ $0.692 \pm 0.008$	$0.704 \pm 0.017$ $0.798 \pm 0.013$ $0.821 \pm 0.003$	$0.236 \pm 0.025$ $0.306 \pm 0.003$ $0.293 \pm 0.003$	$0.470 \pm 0.02$ $0.552 \pm 0.00$ $0.557 \pm 0.00$
	Baseline			0.524	0.476	0.556	0.463	0.505	0.123	0.314
	250	Merging	M-SMoE HC-SMoE	$\begin{array}{c} 0.315 \pm 0.007 \\ 0.439 \pm 0.028 \end{array}$	$\begin{array}{c} 0.270 \pm 0.015 \\ 0.386 \pm 0.020 \end{array}$	$\begin{array}{c} 0.446 \pm 0.007 \\ 0.530 \pm 0.022 \end{array}$	$\begin{array}{c} 0.380 \pm 0.015 \\ 0.441 \pm 0.007 \end{array}$	$\begin{array}{c} 0.353 \pm 0.008 \\ 0.449 \pm 0.005 \end{array}$	$\begin{array}{c} 0.033 \pm 0.010 \\ 0.110 \pm 0.010 \end{array}$	$0.193 \pm 0.00$ $0.279 \pm 0.00$
Mixtral-8x7B- Instruct-v0.1	25%	Pruning	Frequency EAN REAP	$\begin{array}{c} 0.400 \pm 0.034 \\ 0.413 \pm 0.027 \\ 0.439 \pm 0.018 \end{array}$	$\begin{array}{c} 0.358 \pm 0.035 \\ 0.366 \pm 0.024 \\ 0.370 \pm 0.007 \end{array}$	$\begin{array}{c} 0.541 \pm 0.006 \\ 0.477 \pm 0.009 \\ 0.535 \pm 0.011 \end{array}$	$\begin{array}{c} 0.453 \pm 0.012 \\ 0.409 \pm 0.013 \\ 0.452 \pm 0.011 \end{array}$	$\begin{array}{c} 0.438 \pm 0.018 \\ 0.416 \pm 0.015 \\ 0.449 \pm 0.002 \end{array}$	$\begin{array}{c} 0.099 \pm 0.014 \\ 0.111 \pm 0.006 \\ 0.102 \pm 0.010 \end{array}$	$0.269 \pm 0.00$ $0.264 \pm 0.00$ $0.275 \pm 0.00$
	50%	Merging	M-SMoE HC-SMoE	$\begin{array}{c} 0.085 \pm 0.026 \\ 0.175 \pm 0.015 \end{array}$	$\begin{array}{c} 0.076 \pm 0.022 \\ 0.146 \pm 0.000 \end{array}$	$\begin{array}{c} 0.139 \pm 0.121 \\ 0.335 \pm 0.026 \end{array}$	$\begin{array}{c} 0.118 \pm 0.102 \\ 0.282 \pm 0.031 \end{array}$	$\begin{array}{c} 0.091 \pm 0.079 \\ 0.235 \pm 0.018 \end{array}$	$\begin{array}{c} 0.004 \pm 0.006 \\ 0.013 \pm 0.008 \end{array}$	$0.047 \pm 0.03$ $0.124 \pm 0.00$
	30%	Pruning	Frequency EAN REAP	$0.187 \pm 0.015$ $0.220 \pm 0.006$ $0.232 \pm 0.018$	$0.148 \pm 0.007$ $0.189 \pm 0.006$ $0.193 \pm 0.013$	$0.342 \pm 0.016$ $0.375 \pm 0.020$ $0.274 \pm 0.106$	$0.287 \pm 0.012$ $0.325 \pm 0.015$ $0.241 \pm 0.087$	$0.241 \pm 0.007$ $0.277 \pm 0.005$ $0.235 \pm 0.056$	$0.023 \pm 0.004$ $0.031 \pm 0.011$ $0.035 \pm 0.003$	$0.132 \pm 0.00$ $0.154 \pm 0.00$ $0.135 \pm 0.02$
	Baseline			0.829	0.768	0.788	0.640	0.757	0.341	0.549
	25%	Merging	M-SMoE HC-SMoE	0.823 0.787	0.762 0.738	0.786 0.735	0.635 0.587	0.752 0.712	0.324 0.148	0.538 0.430
Llama-4-Scout- 17B-16E- Instruct		Pruning	Frequency EAN REAP	0.835 0.823 0.829	0.768 0.762 0.787	0.788 0.804 0.788	0.630 0.648 0.622	0.755 0.759 0.756	0.317 0.328 0.242	0.536 <b>0.544</b> 0.499
		Merging	M-SMoE HC-SMoE	0.787 0.604	0.732 0.530	0.762 0.500	0.614 0.399	0.723 0.508	0.187 0.077	0.455 0.293
	50%	Pruning	Frequency EAN REAP	0.823 0.805 0.841	0.756 0.744 0.768	0.751 0.754 0.762	0.595 0.601 0.624	0.731 0.726 0.749	0.223 0.209 0.248	0.477 0.468 <b>0.499</b>
	Baseline			0.848	0.829	0.860	0.743	0.820	0.83 ± 0.066 0.112 ± 0.064 0.112 ± 0.064 0.137 ± 0.119 0.302 0.293 ± 0.017 0.258 ± 0.000 0.302 ± 0.011 0.311 ± 0.018 0.308 ± 0.015 0.205 ± 0.019 0.185 ± 0.018 0.235 ± 0.003 0.123 0.033 ± 0.001 0.110 ± 0.010 0.099 ± 0.014 0.111 ± 0.006 0.102 ± 0.010 0.004 ± 0.006 0.103 ± 0.004 0.013 ± 0.004 0.013 ± 0.004 0.013 ± 0.004 0.013 ± 0.004 0.013 ± 0.004 0.014 ± 0.006 0.013 ± 0.004 0.013 ± 0.004 0.013 ± 0.004 0.014 ± 0.006 0.013 ± 0.004 0.014 ± 0.006 0.013 ± 0.004 0.014 ± 0.006 0.013 ± 0.004 0.014 ± 0.006 0.013 ± 0.004 0.014 ± 0.006 0.013 ± 0.004 0.014 ± 0.006 0.015 ± 0.006 0.017 ± 0.006 0.017 ± 0.006 0.017 ± 0.006 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009 0.009	0.597
		Merging	M-SMoE HC-SMoE	0.866 0.872	0.793 0.805	0.807 0.825	0.659 0.669	0.781 0.793		0.555 0.578
GLM-4.5-Air	25%	Pruning	Frequency EAN REAP	0.848 0.872 0.866	0.811 0.817 0.805	0.854 0.876 0.828	0.706 0.720 0.677	0.805 0.821 0.794	0.374	0.573 <b>0.597</b> <u>0.592</u>
	500	Merging	M-SMoE HC-SMoE	0.518 0.707	0.500 0.659	0.519 0.706	0.437 0.577	0.493 0.662		0.296 0.441
	50%	Pruning	Frequency EAN REAP	0.628 0.841 0.878	0.573 0.780 0.841	0.534 0.807 0.712	0.450 0.661 0.587	0.546 0.773 0.755	0.253	0.325 0.513 <b>0.553</b>
	Baseline			0.951	0.890	0.923	0.791	0.889	$0.431 \pm 0.011$	0.660
Qwen3-Coder- 480B-A35B- Instruct-FP8	25%	Pruning	Frequency EAN REAP	0.884 0.939 0.957	0.805 0.878 0.890	0.810 0.911 0.917	0.669 0.775 0.772	0.792 0.876 0.884	$0.419 \pm 0.015$	0.544 0.647 <b>0.650</b>
	50%	Pruning	Frequency EAN REAP	0.020 0.915 0.939	0.012 0.841 0.872	0.007 0.854 0.910	0.003 0.714 0.772	0.011 0.831 0.873	$0.382 \pm 0.012$	0.011 <u>0.607</u> <b>0.644</b>
	Baseline			0.963	0.921	0.913	0.735	0.883		0.659
Kimi-K2- Instruct-	25%	Pruning	Frequency EAN REAP	0.530 0.909 0.957	0.463 0.860 0.921	0.595 0.857 0.918	0.508 0.698 0.759	0.524 0.831 0.889	0.082 0.379 0.440	0.303 0.605 <b>0.664</b>
W4A16	50%	Pruning	Frequency EAN REAP	0.098 0.866 0.915	0.079 0.811 0.884	0.175 0.780 0.899	0.146 0.632 0.754	0.124 0.772 0.863	0.000 0.253 0.429	0.062 0.513 <b>0.646</b>

Table A6: C4 calibrated results for coding and MC tasks.

Model	Compression	Technique	Method	Eval+	Coding LiveCode	Code Avg	ARC-c	ARC-e	BoolQ	Hellaswag	MC MMLU	OBQA	RTE	WinoG.	MC Avg
ERNIE-4.5-21B- A3B-PT	Baseline			0.861	0.231	0.546	0.564	0.782	0.873	0.813	0.737	0.462	0.812	0.724	0.721
		Merging	M-SMoE HC-SMoE	0.065 0.403	0.016 0.099	0.041 <b>0.251</b>	0.497 0.515	0.729 0.728	0.860 0.860	0.723 0.745	0.602 0.649	0.424 0.428	0.801 0.794	0.699 0.694	0.667 <u>0.677</u>
	25%	Pruning	Frequency EAN REAP	0.274 0.282 0.242	0.000 0.000 0.023	0.137 0.141 0.133	0.515 0.528 0.490	0.735 0.750 0.716	0.841 0.853 0.855	0.719 0.790 0.783	0.588 0.558 0.656	0.382 0.442 0.452	0.791 0.783 0.809	0.683 0.706 0.723	0.657 0.676 <b>0.685</b>
		Merging	M-SMoE HC-SMoE	0.000	0.000	0.000	0.297 0.409	0.460 0.615	0.674 0.666	0.449 0.515	0.312 0.489	0.280 0.290	0.671 0.632	0.575 0.580	0.465 0.524
	50%	Pruning	Frequency EAN REAP	0.000 0.007 0.033	0.000 0.003 0.000	0.000 0.005 0.016	0.393 0.451 0.406	0.625 0.676 0.612	0.717 0.742 0.754	0.569 0.687 0.654	0.496 0.474 0.468	0.324 0.398 0.396	0.758 0.736 0.718	0.656	0.563 <b>0.607</b> <u>0.583</u>
	Baseline			0.859	0.302	0.581	0.563	0.790	0.887	0.778	0.779	0.454	0.816	0.702	0.721
	25%	Merging	M-SMoE HC-SMoE	0.000 0.831	0.000 0.269	0.000 <b>0.550</b>	0.551 0.470	0.768 0.713	0.883 0.833	0.761 0.622	0.733 0.646	0.418 0.376	0.848 0.805	0.701 0.665	0.708 0.641
Qwen3-30B-A3B		Pruning	Frequency EAN REAP	0.000 0.000 0.735	0.000 0.000 0.227	0.000 0.000 <u>0.481</u>	0.548 0.569 0.557	0.789 0.802 0.781	0.889 0.889 0.872	0.775 0.774 0.746	0.735 0.735 0.718	0.438 0.438 0.436	0.801 0.801 0.794	0.694 0.697 0.704	0.709 <b>0.713</b> 0.701
	50%	Merging	M-SMoE HC-SMoE	0.000	0.000 0.209	0.000 <b>0.468</b>	0.262 0.316	0.348 0.495	0.693 0.715	0.479 0.354	0.237 0.422	0.290 0.282	0.523 0.603	0.542 0.536	0.422 0.465
		Pruning	Frequency EAN REAP	0.000 0.000 0.006	0.000 0.000 0.000	0.000 0.000 <u>0.003</u>	0.349 0.480 0.421	0.488 0.736 0.640	0.782 0.876 0.837	0.672 0.760 0.653	0.503 0.607 0.495	0.364 0.424 0.388	0.588 0.762 0.704	0.619 0.694 0.635	0.545 <b>0.667</b> 0.596
	Baseline			0.505	0.123	0.314	0.650	0.842	0.887	0.861	0.691	0.496	0.722	0.740	0.736
		Merging	M-SMoE HC-SMoE	0.320 0.420	0.044 0.121	0.182 <b>0.271</b>	0.532 0.608	0.775 0.811	0.828 0.876	0.746 0.838	0.529 0.631	0.424 0.484	0.603 0.736	0.632 0.726	0.634 0.714
Mixtral-8x7B- Instruct-v0.1	25%	Pruning	Frequency EAN REAP	0.396 0.399 0.415	0.070 0.092 0.077	0.233 0.246 0.246	0.612 0.613 0.606	0.816 0.814 0.807	0.868 0.875 0.875	0.836 0.842 0.835	0.593 0.613 0.633	0.482 0.498 0.486	0.675 0.690 0.791	0.739 0.733 0.709	0.703 0.710 <b>0.718</b>
		Merging	M-SMoE HC-SMoE	0.000	0.000 0.033	0.000 <b>0.103</b>	0.260 0.540	0.460 0.764	0.614 0.862	0.395 0.795	0.240 0.544	0.302 0.448	0.527 0.675	0.526 0.709	0.416 0.667
	50%	Pruning	Frequency EAN REAP	0.173 0.139 0.167	0.008 0.008 0.012	0.090 0.074 0.089	0.504 0.550 0.525	0.739 0.756 0.774	0.793 0.842 0.856	0.771 0.804 0.794	0.463 0.529 0.533	0.426 0.460 0.454	0.675 0.726 0.751	0.646 0.716 0.688	0.627 <b>0.673</b> 0.672

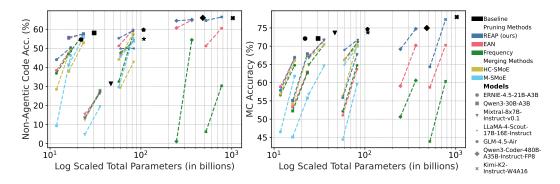


Figure A5: Coding and MC accuracy across all models vs. parameters. The benefits of REAP over other compression methods are evident at 50% compression. For large-scale SMoEs, REAP is near-lossless whereas the shortcomings of frequency-based pruning become apparent.

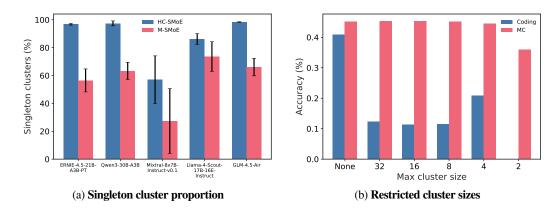


Figure A6: (a) Average proportion of singleton clusters vs. model for HC-SMoE and M-SMoE. We find that the clustering algorithms used by our baseline merging methods tend to generate a high proportion of singleton clusters containing just a single expert. In order to achieve the desired compression ratio, the large number of singletons conversely results in some clusters which contain many experts, in some cases N/2+1 experts for a layer with N experts are grouped into a single cluster. (b) Accuracy vs. maximum cluster size using M-SMoE to compress 50% of experts in Qwen3-30B. While MC accuracy remains stable up to a maximum cluster size of 4, generative coding capabilities are severely diminished by restricting the clustering algorithm.

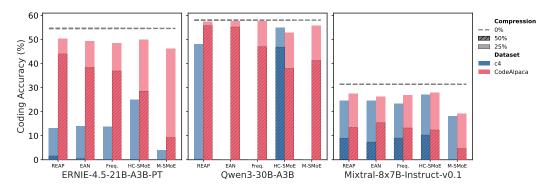


Figure A7: **Coding accuracy vs. calibration dataset**. Using domain-specific calibration datasets substantially improves compressed model quality within the target domain. Fine-grained models such as Qwen3-30B and ERNIE suffers greater degradation, with several compression methods failing to produce any coherent output when calibrated on c4.

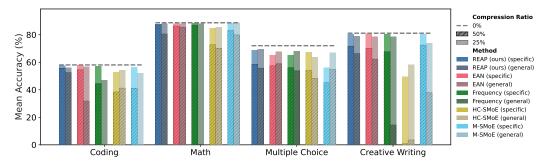


Figure A8: Mean accuracy vs. task type for models calibrated with domain specific data versus general data. The "general" calibration data consists of the combination of evol-codealpaca-v1, Writing-Prompts curated, and tulu-3-sft-personas-math and includes three times the total number of samples as the domain-specific calibration datasets. While the general data calibrated models perform reasonably well at 25% compression, domain-specific data is crucial for high-quality compressed SMoE accuracy at 50% compression.