

# 000 001 SCALING LAWS MEET MODEL ARCHITECTURE: TO- 002 WARD INFERENCE-EFFICIENT LLMs 003 004

005 **Anonymous authors**

006 Paper under double-blind review

## 007 008 ABSTRACT 009

011 Scaling the number of parameters and the size of training data has proven to be  
012 an effective strategy for improving large language model (LLM) performance.  
013 Yet, as these models grow increasingly powerful and widely deployed, the cost  
014 of inference has become a pressing concern. Despite its importance, the trade-  
015 off between model accuracy and inference efficiency remains underexplored. In  
016 this work, we examine how key architectural factors, hidden size, the allocation  
017 of parameters between MLP and attention (mlp-to-attention ratio), and grouped-  
018 query attention (GQA), influence both inference cost and accuracy. We introduce a  
019 conditional scaling law that augments the Chinchilla framework with architectural  
020 information, along with a search framework for identifying architectures that are  
021 simultaneously inference-efficient and accurate. To validate our approach, we  
022 train more than 200 models spanning 80M to 3B parameters and 8B to 100B  
023 training tokens, and fit the proposed conditional scaling law. Our results show  
024 that the conditional scaling law reliably predicts optimal architectural choices and  
025 that the resulting models outperform existing open-source baselines. Under the  
026 same training budget, optimized architectures achieve up to 2.1% higher accuracy  
027 and 42% greater inference throughput compared to LLaMA-3.2.

## 028 1 INTRODUCTION 029

030 Scaling law studies Kaplan et al. (2020); Hoffmann et al. (2022); Muennighoff et al. (2023); Krajew-  
031 ski et al. (2024); Abnar et al. (2025) have shown that increasing model parameters, training tokens,  
032 dataset quality, and compute budget consistently reduces pre-training loss, improves downstream  
033 task performance Hendrycks et al. (2021); Austin et al. (2021), and enables the emergence of novel  
034 capabilities Wei et al. (2022). These insights have driven the development of many state-of-the-art  
035 large language models Touvron et al. (2023); Yang et al. (2025); Guo et al. (2025).

036 However, as the field advances, it has become increasingly clear that focusing exclusively on training  
037 overlooks the practical challenges of deploying these models at scale Chien et al. (2023); Wu  
038 et al. (2024); Muhamed et al. (2023). A major limitation of existing scaling laws is their omission  
039 of inference costs, which constitute the dominant expense in deploying large models in real-world  
040 applications Sardana et al. (2023); Park et al. (2024). Moreover, the growing use of LLMs in reasoning  
041 systems highlights the need for scaling laws that account for inference costs Snell et al. (2024);  
042 Brown et al. (2024); Luo et al. (2024); Qi et al. (2024); Guan et al. (2025). Therefore, we ask the  
043 following question:

044 *045 Can we explicitly capture the trade-off between inference efficiency and accuracy  
046 of large language models?*

047 To address this question, a recent study Sardana et al. (2023) proposed scaling laws that incorporate  
048 the total FLOPs from both training and inference. However, their formulation requires estimating  
049 the total number of tokens generated over a model’s entire lifespan. Because inference is performed  
050 repeatedly during deployment, this assumption renders the proposed scaling law impractical for  
051 real-world use. Another study Bian et al. (2025) extends Chinchilla scaling laws by incorporating  
052 model architecture. However, this work has notable limitations. First, the study considers only  
053 the aspect ratio, defined as hidden size over number of layers, as the architectural factor. Yet, as  
shown in Figure 1, aspect ratio alone fails to capture the full range of factors that influence inference

efficiency in large language models. Second, the depth of the model strongly influences accuracy: cutting layers tends to impair the model’s generalization after fine-tuning Petty et al. (2023). Finally, the study lacks a general framework for incorporating broader architectural factors, including hidden size and GQA, into scaling laws.

In this work, we fix the number of layers and study the effect of other architectural factors, including GQA, hidden size, and the mlp-to-attention ratio. This design choice is motivated by recent open-weight models such as LLaMA Touvron et al. (2023), Qwen Yang et al. (2025), Gemma Team et al. (2024a), and Phi Abdin et al. (2024), which, despite having a comparable number of parameters, adopt markedly different architectural designs.

Our primary goal is to investigate how model architecture influences both inference efficiency and model accuracy. We begin by comparing the inference efficiency of models with identical parameter counts but varying architectures. Next, we train over 200 models, ranging from 80M to 297M parameters on up to 30B tokens, to systematically characterize the relationship between architectural design and accuracy. Guided by these empirical findings, we introduce a conditional extension of the Chinchilla scaling laws that incorporates architectural parameters, establishing a general framework for identifying model architectures that balance inference efficiency and performance.

Finally, we validate this framework by fitting the proposed scaling law on models between 80M and 297M parameters, and evaluating its predictions when scaling up to pretrain 3B-parameter models. Our results demonstrate that, under identical training setups, the derived optimal 3B-parameter architecture achieves up to 42% higher inference throughput than the LLaMA-3.2-3B architecture, while maintaining better accuracy.

## 2 BACKGROUND

Accurately predicting the performance of large language models during scaling is essential. This enables us to answer key questions: (i) what is the optimal allocation of available resources between model size and training tokens, and (ii) what performance gains can be expected from additional resources? Fortunately, the model loss has been observed to follow a power-law relationship with respect to the number of parameters  $N$  and training tokens  $D$  Hoffmann et al. (2022); Muennighoff et al. (2023) with:

$$L(N, D) = E + \frac{A}{N^\alpha} + \frac{B}{D^\beta} \quad (1)$$

where  $L$  is the model loss,  $N$  is the number of total parameters and  $D$  is the number of tokens used for training and  $A, B, E, \alpha, \beta$  are parameters to be learned.

To fit the learnable parameters in Eq. (1), Chinchilla Hoffmann et al. (2022) employs two strategies: (i) training models with a fixed number of parameters while varying the number of training tokens, and (ii) training models under a fixed compute budget<sup>1</sup>, varying both parameters and tokens. The resulting data are combined to fit the learned parameters in Eq. (1). With the fitted scaling laws, Chinchilla addresses the following question to determine optimal allocation:

$$\arg \min_{N, D} L(N, D) \text{ s.t. } \text{FLOPs}(N, D) = C \quad (2)$$

where  $C$  denotes the resource constraint,  $N$  the total number of parameters, and  $D$  the number of training tokens.

<sup>1</sup>The compute cost is approximated as  $\text{FLOPs}(N, D) \approx 6ND$  in Hoffmann et al. (2022); Muennighoff et al. (2023), where  $N$  denotes the number of parameters and  $D$  the number of training tokens. In this work, we adopt the same settings as prior studies.

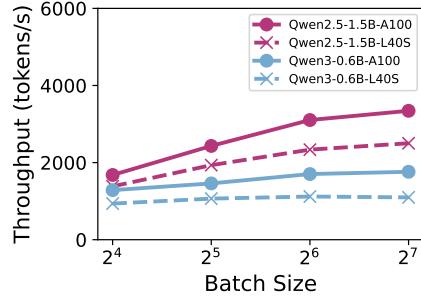


Figure 1: Although larger models generally achieve lower inference throughput than smaller ones, Qwen2.5-1.5B outperforms Qwen3-0.6B. Despite having the same number of layers, Qwen2.5-1.5B benefits from a higher hidden size, GQA, and mlp-to-attention ratio.

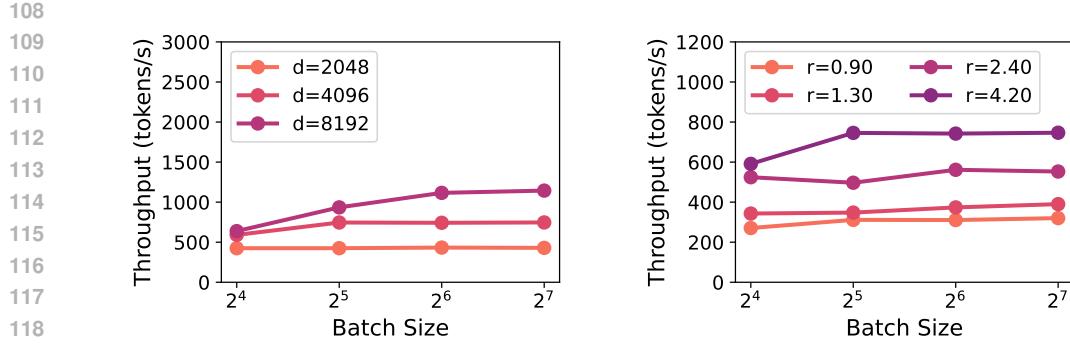


Figure 2: **Inference throughput.** (left) hidden size  $d = d_{\text{model}}$  and (right) mlp-to-attention ratio  $r = r_{\text{mlp}/\text{attn}}$  on the 8B model. Under a fixed parameter budget  $N_{\text{non-embed}}$ , larger hidden sizes and higher mlp-to-attention ratios improve inference throughput for varying batch sizes.

In this paper, we do not address how to optimally allocate compute between model size and training data under a fixed compute budget. Instead, our focus is on identifying model architectures that optimize inference efficiency and accuracy under fixed parameter and token budgets. For example, given a model with 7B parameters trained on 14T tokens, we study how to design an architecture that satisfies both efficiency and accuracy requirements.

### 3 MODEL ARCHITECTURE-AWARE SCALING LAWS

#### 3.1 MODEL ARCHITECTURE VARIATIONS

The architecture of a decoder-only transformer is composed of a sequence of stacked decoder blocks, each sharing the same structure to facilitate model-parallel deployment across devices. Under this design, the overall architecture of dense LLMs is primarily determined by the hidden size and the MLP intermediate size, which together specify the attention and MLP layers structure. This work studies the optimal model architecture given a fixed total number of non-embedding parameters  $N_{\text{non-embed}}$  (at different levels). Although the number of layers  $n_{\text{layer}}$  also plays a critical role (closely related to aspect ratio (Petty et al., 2023)), varying  $n_{\text{layer}}$  under a fixed  $N_{\text{non-embed}}$  substantially impacts both inference cost and accuracy (Tay et al., 2021; Alabdulmohsin et al., 2023). Therefore, we fix  $n_{\text{layer}}$  and focus on the effects of hidden size  $d_{\text{model}}$  and the mlp-to-attention ratio  $r_{\text{mlp}/\text{attn}}$  on inference efficiency (§3.2) and accuracy (§3.3), noting that  $n_{\text{layer}}$  still varies across different  $N_{\text{non-embed}}$  levels. In §3.3, we introduce a conditional scaling law to predict the performance of architectural variants, and in §3.4, we present a lightweight framework for identifying architectures that optimally balance inference efficiency and accuracy.

Note that the number of attention parameters is primarily determined by the hidden size  $d_{\text{model}}$  and the attention projection dimension, since most open-weight models adopt non-square  $q, k, v$  projection matrices, as seen in Gemma (Team et al., 2024a) and Qwen3 (Yang et al., 2025). For consistency, we fix the per-head dimension  $d_{\text{head}}$  to 64 for models with  $N_{\text{non-embed}} \leq 1\text{B}$  and to 128 for models with  $N_{\text{non-embed}} \geq 3\text{B}$ . Consequently, to maintain a constant  $r_{\text{mlp}/\text{attn}}$ , we adjust the number of attention heads  $n_{\text{head}}$  rather than altering the projection dimension directly. This design choice also provides flexibility to incorporate architectural variants such as grouped-query attention.

#### 3.2 INFERENCE EFFICIENCY

Inspired by the success and widespread adoption of open-weight dense models such as Qwen3 (Yang et al., 2025), LLaMA-3.2 (Dubey et al., 2024), and the Gemma-2 (Team et al., 2024b) family, we construct architectural variants by modifying the configurations of the LLaMA-3.2 and Qwen3 dense models (Figure 11-13 in Appendix E). In addition to hidden size and the mlp-to-attention ratio, we find that group-query attention has a critical impact on inference efficiency, even though it only modestly reduces the number of attention parameters (by shrinking the key and value matrices). To

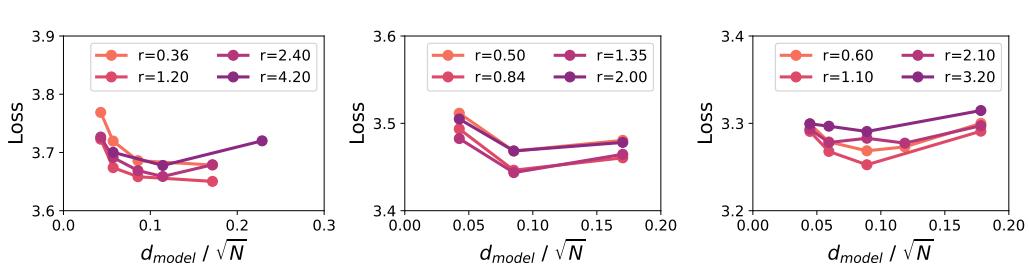


Figure 3: **Loss vs. hidden size.** (Left) 80M model variants; (Center) 145M model variants; (Right) 297M model variants. Across model sizes, the relationship between training loss and  $d_{\text{model}} / \sqrt{N}$  exhibits a consistent U-shaped curve when architectural factors such as GQA and the MLP-to-attention ratio are held fixed. The legend denotes the MLP-to-attention ratio  $r = r_{\text{mlp}/\text{attn}}$  for each model.

disentangle these effects, we conduct controlled ablations of hidden size, MLP-to-attention ratio, and GQA under the following setups:

- *hidden size  $d_{\text{model}}$ :* fix  $N_{\text{non-embed}}$ ,  $r_{\text{mlp}/\text{attn}}$  and GQA = 4, vary  $d_{\text{model}}$  and number of attention heads  $n_{\text{head}}$  (Figure 2 left).
- *mlp-to-attention ratio  $r_{\text{mlp}/\text{attn}}$ :* fix  $N_{\text{non-embed}}$ ,  $d_{\text{model}}$  and GQA = 4, vary  $n_{\text{head}}$  and intermediate size (Figure 2 right).
- *GQA:* fix  $N_{\text{non-embed}}$ ,  $d_{\text{model}}$  and  $r_{\text{mlp}/\text{attn}}$ , vary  $n_{\text{head}}$  and number of key-value heads (Appendix E).

Figure 2 shows the ablation of varying hidden sizes  $d_{\text{model}}$  and mlp-to-attention  $r_{\text{mlp}/\text{attn}}$  on the LLaMA-3.1-8B model variants. We observe that larger hidden size (or fewer attention heads) and higher mlp-to-attention ratios improve inference throughput. Similar trends are observed in the LLaMA-3.2-1B and 3B model variants (Appendix E). These gains arise in part because larger  $d_{\text{model}}$  and higher  $r_{\text{mlp}/\text{attn}}$  reduce the total FLOPs, as detailed in the inference FLOPs analysis (Appendix J). In addition, these architectural choices shrink the KV cache, lowering I/O cost during inference and further improving throughput Adnan et al. (2024). Figure 10 in Appendix E presents the GQA ablation, confirming prior observations Ainslie et al. (2023) that increasing GQA consistently improves inference throughput. A comparable set of ablation experiments on Qwen3 models, also reported in Appendix E, further corroborates these findings.

### 3.3 A CONDITIONAL SCALING LAW

Improving inference efficiency should not come at the expense of significantly reducing model accuracy, making it crucial to understand how architectural choices affect accuracy and training loss. Because training large-scale language models is prohibitively expensive, a common strategy is to study smaller models and use scaling laws to extrapolate insights to larger scales, for example, the Chinchilla scaling laws (Hoffmann et al., 2022). However, incorporating multiple architectural factors into such laws remains challenging. To address this, we examine the effect of architectural choices on training loss  $L$  in a conditional manner, varying one factor at a time while keeping the others fixed.

**hidden size  $d_{\text{model}}$ .** We note that  $d_{\text{model}}$  generally scales linearly with  $\sqrt{N_{\text{non-embed}}}$ . Assuming squared attention weight matrices, the number of attention parameters  $N_{\text{attn}}$  can be expressed as

$$4d_{\text{model}}^2 \propto N_{\text{attn}} = N_{\text{non-embed}} \times \frac{r}{r+1},$$

where  $r = r_{\text{mlp}/\text{attn}}$  is fixed, and the constant factor 4 arises from the query, key, value, and output projection layers in each attention block. To capture this scaling behavior, we normalize  $d_{\text{model}}$  by  $\sqrt{N_{\text{non-embed}}}$  and examine its relation to loss  $L$  in Figure 3. The resulting U-shaped curves  $L(d/\sqrt{N} | r, N, D)$  exhibit nearly identical optima across different model sizes. Moreover, Figure 3 confirms that excessively large hidden sizes, which reduce the number of attention heads  $n_{\text{head}}$ , can degrade accuracy—a phenomenon consistently observed in prior analyses of transformer capacity and head allocation (Kaplan et al., 2020; Hoffmann et al., 2022).

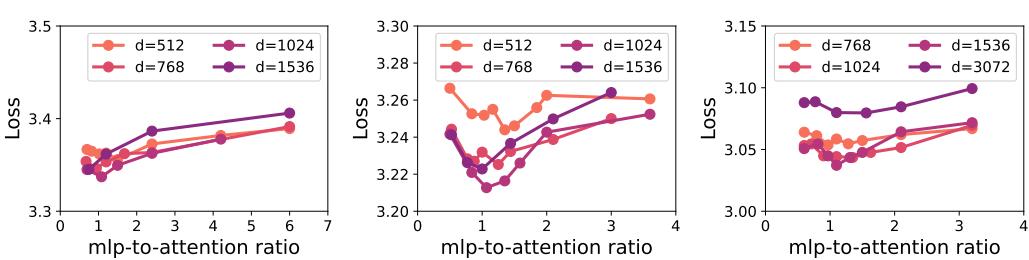


Figure 4: **Loss vs. MLP-to-attention ratio.** (Left) 80M model variants; (Center) 145M model variants; (Right) 297M model variants. Across model sizes, the relationship between training loss and  $r_{\text{mlp}/\text{attn}}$  exhibits a consistent U-shaped curve when architectural factors such as GQA and hidden size are held fixed. The legend denotes the hidden size  $d = d_{\text{model}}$  for each model.

**mlp-to-attention ratio  $r_{\text{mlp}/\text{attn}}$ .** Figure 4 illustrates how the loss varies with  $r_{\text{mlp}/\text{attn}}$ , conditioned on  $d_{\text{model}}$  fixed at different levels, where we consistently observe a U-shaped curve  $L(r | d/\sqrt{N}, N, D)$ . While the attention mechanism is central to the success of transformers (Vaswani, 2017), recent open-weight models have allocated a progressively smaller fraction of parameters to attention as overall model size increases (e.g., LLaMA and Qwen families). Our analysis indicates that this trend is not universally optimal: there exists an interior optimum in the allocation of attention parameters, and deviating from it in either direction degrades model performance. This suggests that careful tuning of the mlp-to-attention ratio is critical for scaling transformers effectively.

As shown in Figures 3 and 4, both hidden size and the MLP-to-attention ratio exhibit U-shaped relationships with training loss. To capture these trends, we fit the function  $c_0 + c_1 \log x + c_2/x$  separately for  $x = r_{\text{mlp}/\text{attn}}$  and  $d_{\text{model}}/\sqrt{N_{\text{non-embed}}}$ . This formulation effectively models the U-shaped behavior while ensuring sublinear growth as  $x$  increases. However, incorporating  $r_{\text{mlp}/\text{attn}}$ ,  $d_{\text{model}}$ ,  $N$ , and  $D$  into a unified, architecture-aware scaling law remains challenging. Since fitting a single all-purpose scaling law  $L(d/\sqrt{N}, r, N, D)$  is unrealistic across all possible configurations, we instead propose a two-step conditional approach:

1. For given  $N$  and  $D$ , obtain the optimal loss  $L_{\text{opt}}(N, D) = \min L(N, D) = \min (E + \frac{A}{N^\alpha} + \frac{B}{D^\beta})$  from the Chinchilla scaling law (Eq. 1) as a reference point.
2. Calibrate the loss of architectural variants  $L(d/\sqrt{N}, r | N, D)$  relative to this reference.

We focus on two simple and transparent calibration schemes:

- (multiplicative)

$$L(d/\sqrt{N}, r | N, D) = (a_0 + a_1 \log(\frac{d}{\sqrt{N}}) + a_2 \frac{\sqrt{N}}{d}) \cdot (b_0 + b_1 \log r + \frac{b_2}{r}) \cdot L_{\text{opt}} \quad (3)$$

- (additive)  $L(d/\sqrt{N}, r | N, D) = (a_0 + a_1 \log(\frac{d}{\sqrt{N}}) + a_2 \frac{\sqrt{N}}{d}) + (b_1 \log r + \frac{b_2}{r}) + L_{\text{opt}}$

Here,  $a_i$  and  $b_i$  are learnable parameters that are shared across all  $N, D$ . Note that both functional forms assume the effects of  $r_{\text{mlp}/\text{attn}}$  and  $d_{\text{model}}$  on loss are separable.

### 3.4 SEARCHING FOR INFERENCE-EFFICIENT ACCURATE MODELS

With the conditional scaling law, we can identify architectures that are both inference-efficient and accurate by solving the following optimization problem: given  $N, D$ , and a set of architectural choices  $P$ ,

$$\text{argmax}_P I_N(P), \quad \text{s.t.} \quad L(P | N, D) \leq L_t, \quad (4)$$

where  $I_N(P)$  denotes the inference efficiency of an architecture  $P$  with total  $N_{\text{non-embed}}$  parameters, and  $L_t$ , ( $\geq L_{\text{opt}}$ ) is the maximum allowable training loss.

As shown in Figure 10 (Appendix E), GQA has a substantial impact on inference efficiency; However, unlike hidden size and the mlp-to-attention ratio, GQA does not exhibit a consistent continuous relationship with loss (Figure 23, Appendix H) and is highly variable, making it challenging to identify settings that achieve both accuracy and efficiency. Fortunately, the search space for GQA is

270 relatively small once  $N_{\text{non-embed}}$ ,  $d_{\text{model}}$ , and  $r_{\text{mlp/attn}}$  are fixed, since GQA must be a prime factor of  
 271 the number of attention heads  $n_{\text{head}}$ . In practice, we perform a local GQA search by enumerating  
 272 feasible values and applying early stopping once performance falls below that of the GQA= 4 base-  
 273 line. Algorithm 1 summarizes our overall framework for identifying inference-efficient and accurate  
 274 architectures.

275

---

**276 Algorithm 1:** Searching for Inference-Efficient Accurate Model
 

---

277 **Input:** Model parameters  $N$ , training tokens  $D$ , target loss  $L_t$ ; inference efficiency  $I_N(\cdot)$ ;  
 278 optional: the optimal loss  $L_{\text{opt}}(N, D)$

279 Train smaller models to fit the Chinchilla scaling laws (Eq. 1) if  $L_{\text{opt}}(N, D)$  is unavailable

280 Solve the constrained optimization (Eq. 4) for  $d_{\text{model}}$ ,  $r_{\text{mlp/attn}}$  and corresponding architecture  $P$

281 Perform a local search over GQA values with early stopping to maximize inference efficiency

282 **return** Final model architecture  $\{P, \text{GQA}\}$

---

283

284

285 **4 EXPERIMENT SETUP**

286

287 We first detail the experimental setup of training, inference, and downstream task evaluation, and  
 288 then describe how we derive the conditional scaling law and scale up to larger sizes.

289

290 **Training Setup.** We sample the training data from Dolma-v1.7 Soldaini et al. (2024), which  
 291 contains data from 15 different sources. Tokens are sampled with probability proportional to  
 292 each source’s contribution, ensuring the sampled dataset preserves a similar distribution to Dolma-  
 293 v1.7. We train decoder-only LLaMA-3.2 (Dubey et al., 2024) style transformers with  $N_{\text{non-embed}}$  in  
 294  $\{80\text{M}, 145\text{M}, 297\text{M}, 1\text{B}, 3\text{B}\}$ , for each  $N_{\text{non-embed}}$ , we obtain model architecture candidates by varying  
 295 hidden size  $d_{\text{model}}/\sqrt{N_{\text{non-embed}}}$  and mlp-to-attention ratio  $r_{\text{mlp/attn}}$ . (changing intermediate size  
 296 and number of attention heads  $n_{\text{head}}$ ) while holding other architectural factors fixed e.g. GQA= 4.  
 297 A full list of over 200 model architectures used can be found in Appendix C. All models are trained  
 298 on  $100N_{\text{non-emb}}$  tokens ( $5 \times$  Chinchilla optimal) to ensure convergence. We tuned training hyper-  
 299 parameters (mainly following prior work Chen et al. (2025)), with a full list in Appendix D.

300

301 **Inference Setup.** We evaluate the inference efficiency using the vLLM framework Kwon et al.  
 302 (2023). By default, inputs consist of 4096 tokens and outputs of 1024 tokens. We report the av-  
 303 eraged inference throughput (tokens/second) from 5 repeated runs. Unless otherwise specified, all  
 304 experiments are conducted on NVIDIA Ampere A100 GPUs (40GB) with vLLM.

305

306 **LLM Evaluation Setup.** Following prior works Biderman et al. (2023); Zhang et al. (2024),  
 307 we evaluate pretrained models in the zero-shot setting using lm-evaluation-harness<sup>2</sup> on  
 308 nine benchmarks: ARC-Easy Clark et al. (2018), ARC-Challenge Clark et al. (2018), LAM-  
 309 BADA Paperno et al. (2016), HellaSwag Zellers et al. (2019), OpenBookQA Mihaylov et al. (2018),  
 310 PIQA Bisk et al. (2020), SciQ Welbl et al. (2017), WinoGrande Sakaguchi et al. (2021), and  
 311 CoQA Reddy et al. (2019).

312

313 **Fitting Scaling Laws.** Following Gadre et al. (2024); Bian et al. (2025), we use the Levenberg-  
 314 Marquardt algorithm to fit the conditional scaling laws (Eq. 3). The Levenberg–Marquardt algorithm  
 315 does least-squares curve fitting by estimating  $\hat{\beta}$  as the solution to  $\arg \min_{\beta} \sum_{i=1}^m [y_i - f(x_i, \beta)]^2$ ,  
 316 where  $(x_i, y_i)$  are the observed data pairs. Note that instead of fitting the Chinchilla scaling  
 317 law, we empirically searched over architecture variants to find the optimal loss  $L_{\text{opt}}(N, D)$  for  
 318  $N_{\text{non-embed}} < 1\text{B}$  scale.

319

We scale up the scale law fitting in the following progressive manner:

- 320 (Task 1) fit on the 80M results and evaluate on 145M results;  
 321 (Task 2) fit on 80, 145M results and evaluate on 297M results;  
 322 (Task 3) fit on 80, 145, 297M results and evaluate on 1B results;

323

---

<sup>2</sup><https://github.com/EleutherAI/lm-evaluation-harness>

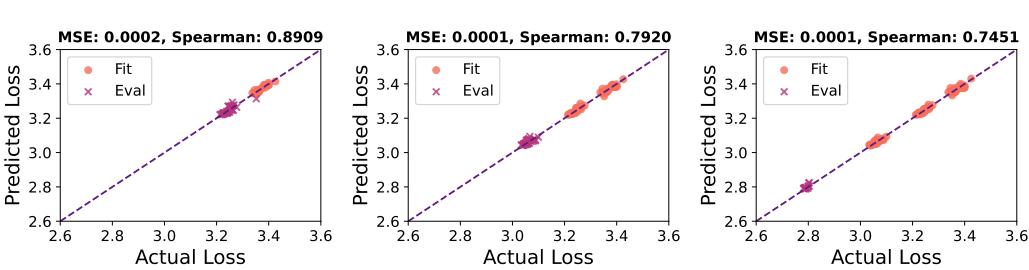


Figure 5: **Predictive performances** of the fitted conditional scaling law on: (left) Task 1: Fit on 80M, evaluate on 145M; (center) Task 2: Fit on 80, 145M, evaluate on 297M; (right) Task 3: Fit on 80, 145, 297M, evaluate on 1B. Orange dots denote fitting data points, and purple crosses indicate the test data points. We compare scaling-law predicted loss with actual pretraining loss of architectures and observed a consistently low MSE and high Spearman correlation across model scales.

This ensures a robust and consistent way of scaling up the model sizes and evaluating our conditional scaling law. Following prior work Kumar et al. (2024), we evaluate the fitted scaling law with mean squared error (MSE) metric, defined as  $\frac{1}{n} \sum_{i=1}^n (l_i - \hat{l}_i)^2$  where  $l_i$  denotes the actual loss and  $\hat{l}_i$  the predicted loss. We additionally report the Spearman’s rank correlation coefficient Spearman (1961) to compare predicted and actual rankings. Both metrics are calculated on the val data points.

## 5 EXPERIMENT RESULTS

We begin by evaluating the predictive performances of the conditional scaling laws with multiplicative calibration. We then conduct ablation studies to assess the impact of data selection and to evaluate the performance of the scaling laws under additive calibration. Finally, we apply the fitted scaling laws to guide the training of large-scale models following the search framework (§5.1).

**Predictive Accuracy.** As Task 1-3 described in §4, we fit the conditional scaling laws on 80M, (80M, 145M), and (80M, 145M, 297M) loss-architecture data points, and subsequently evaluate on 145M, 297M, and 1B data, respectively. In Figure 5, the low MSE and high Spearman correlation in tasks across different model scales validate the effectiveness and strong predictive performance of the proposed conditional scaling laws.

**Ablation of Outliers.** The mlp-to-attention ratio  $r_{\text{mlp}/\text{attn}}$  of open-weights models typically fall between 0.5 and 5, for example, the mlp-to-attention ratio for LLaMA-3.2-1B, LLaMA-3.2-3B, and Qwen3-8B are 4.81, 1.5, and 4.67, respectively. In Figure 5, we fit the conditional scaling law using only model architectures with  $r_{\text{mlp}/\text{attn}} \in [0.5, 5]$ . We ablate this choice by training model architectures with outlier  $r_{\text{mlp}/\text{attn}}$  below 0.5 and above 5 (such as 0.1, 12.6) in Appendix C. In Figure 24 (left) and Figure 24 (center) in Appendix I, we show on Task 3 a comparison of fitting the conditional scaling law without and with these outliers (with a clear Spearman correlation score degradation), which suggests to exclude extreme outliers for better predicted performances.

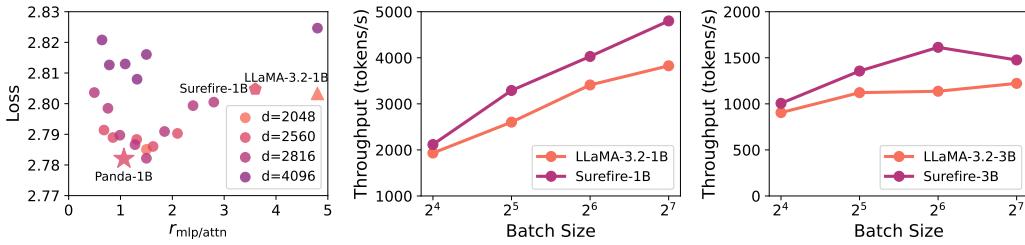
**Ablation of Calibration.** In Figure 24 (right), we ablate an alternative formulation of the scaling laws with additive calibration, as discussed in §3.3. The results on Task 3 show that multiplicative and additive calibrations achieve similar MSE and Spearman correlations. Note that, unlike the conventional unified formulation, both calibrations assume that the effects of  $r_{\text{mlp}/\text{attn}}$  and  $d_{\text{model}}$  on loss are separable. We further ablate more complex joint, non-separable formulations in Appendix I and find that they do not provide superior predictive performance. The two-step reference-and-calibration framework appears robust enough that simple calibrations perform well.

### 5.1 OPTIMAL MODEL ARCHITECTURE

**Validating the conditional scaling law.** We validate the conditional scaling law at the 1B scale by applying multiplicative calibration on Task 3 using data from the (80M, 145M, and 297M) model

378 **Table 1: Large-Scale Model Results.** We evaluate the scaling laws at 1B and 3B scales by training  
 379 Panda-1B, Surefire-1B, and Panda-3B, and compare them with LLaMA-3.2-1B and LLaMA-3.2-  
 380 3B, respectively. The Avg. column reports the mean accuracy across the nine downstream tasks.  
 381 Panda-1B and 3B are trained using the optimal architectural configurations predicted by our scaling  
 382 laws, whereas Surefire-1B and 3B satisfy the loss constraint in Eq. (4) and achieve Pareto optimality.  
 383

| Models       | $d_{\text{model}}$ | $f_{\text{size}}$ | $n_{\text{layers}}$ | GQA | $d_{\text{model}}/\sqrt{N}$ | $r$  | Loss ( $\downarrow$ ) | Avg. ( $\uparrow$ ) |
|--------------|--------------------|-------------------|---------------------|-----|-----------------------------|------|-----------------------|---------------------|
| LLaMA-3.2-1B | 2048               | 8192              | 16                  | 4   | 0.066                       | 4.80 | 2.803                 | 54.9                |
| Panda-1B     | 2560               | 4096              | 16                  | 4   | 0.082                       | 1.07 | 2.782                 | 57.0                |
| Surefire-1B  | 2560               | 6144              | 16                  | 9   | 0.082                       | 3.6  | 2.804                 | 55.4                |
| LLaMA-3.2-3B | 3072               | 8192              | 28                  | 3   | 0.058                       | 4.80 | 2.625                 | 61.9                |
| Panda-3B     | 4096               | 4096              | 28                  | 3   | 0.077                       | 1    | 2.619                 | 62.5                |
| Surefire-3B  | 4096               | 4096              | 28                  | 7   | 0.077                       | 1    | 2.620                 | 62.6                |



402 **Figure 6: Results for 1B and 3B models.** (Left) Panda-1B closely follows the scaling law pre-  
 403 dictions for minimizing training loss. (Center & Right) Inference throughput comparison between  
 404 LLaMA-3.2 and Surefire models, where Surefire is consistently efficient across all batch sizes.  
 405  
 406  
 407

variants. The learned parameters are

$$a_0 = 2.697, a_1 = 0.0974, a_2 = 0.0078, b_0 = 0.3870, b_1 = 0.0063, \text{ and } b_2 = 0.0065.$$

408 From this, we obtain the optimal architectural configuration of  $d_{\text{model}}/\sqrt{N} = 0.08, r = 1.032$  for  
 409 1B model by solving  $\frac{\partial L}{\partial d_{\text{model}}} = 0$  and  $\frac{\partial L}{\partial r} = 0$ . Using this configuration, we train a LLaMA-3.2-style  
 410 1B dense model on 100B tokens, denoted as Panda-1B. Panda-1B outperforms the open-weight  
 411 LLaMA-3.2-1B baseline configs by 2.1% on average across downstream tasks (Table 1). Figure 6  
 412 (left) further confirms the effectiveness of the conditional scaling law by showing that Panda-1B  
 413 achieves the lowest training loss among the exhaustively trained 1B variants under the same setup.  
 414

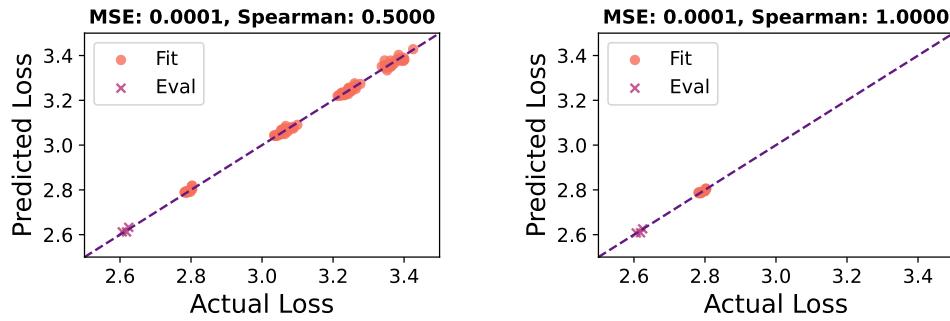
415 We also scale up our methodology to 3B models. Using the same approach but with data from  
 416 the 80M, 145M, 297M, and 1B variants, we fit the scaling law and obtain  $d_{\text{model}}/\sqrt{N} = 0.08$  and  
 417  $r = 1.055$  for the Panda 3B model. Trained on 100B tokens, Panda-3B outperforms the open weight  
 418 LLaMA-3.2-3B configuration by 0.6% on average across downstream tasks (Table 1).  
 419

420 With all components in place, we apply the search framework for inference-efficient and accurate  
 421 models (Alg. 1). For the  $N_{\text{non-embed}} = 1B$  and 3B setting trained on 100B tokens, we set the target  
 422 loss  $L_t$  to match the training loss achieved by the LLaMA-3.2-1B and LLaMA-3.2-3B architectures,  
 423 respectively.

424 **Ablation of inference efficiency.** Although inference efficiency  $I_N(P)$  could, in principle, be  
 425 expressed analytically, it depends heavily on hardware and inference configurations. Therefore,  
 426 rather than solving for  $I_N(P)$  directly, we search over feasible configurations  $P_i$  that satisfy the  
 427 loss constraint on A100 with vLLM and select Pareto-optimal points, which we denote as Surefire-  
 428 1B and Surefire-3B. Surefire-1B and Surefire-3B outperform LLaMA-3.2-1B and LLaMA-3.2-3B  
 429 on downstream tasks (Table 1 with details in Appendix K) and deliver up to 42% higher inference  
 430 throughput (Figure 6, center and right). We also ablate inference efficiency using both vLLM and  
 431 SGLang Zheng et al. (2023) on A100 and NVIDIA H200 GPUs (Appendix E, F). The results remain  
 432 consistent with our vLLM-A100 evaluation: Surefire-1B and 3B outperform LLaMA-3.2-1B and

432 Table 2: **3B Model Ablations.** We assess the robustness of fitting-data strategy at 3B scale by  
 433 training Panda-3B (using 80M, 145M, and 297M data) and Panda-3B<sup>◦</sup> (using only on 1B data), and  
 434 compare both with LLaMA-3.2-3B. Avg. denotes mean accuracy across nine downstream tasks.  
 435

| Models                | $d_{\text{model}}$ | $f_{\text{size}}$ | $n_{\text{layers}}$ | GQA | $d_{\text{model}}/\sqrt{N}$ | $r$  | Loss ( $\downarrow$ ) | Avg. ( $\uparrow$ ) |
|-----------------------|--------------------|-------------------|---------------------|-----|-----------------------------|------|-----------------------|---------------------|
| LLaMA-3.2-3B          | 3072               | 8192              | 28                  | 3   | 0.058                       | 4.80 | 2.625                 | 61.9                |
| Panda-3B              | 4096               | 4096              | 28                  | 3   | 0.077                       | 1    | 2.619                 | 62.5                |
| Panda-3B <sup>◦</sup> | 4096               | 4608              | 28                  | 3   | 0.076                       | 1.23 | 2.606                 | 62.5                |



444 Figure 7: **Effect of the Fitting Data Strategy on Predictive Performance.** (left) Fit on 80M,  
 445 145M, 297M, 1B, evaluate on 3B; (right) Fit on 1B, evaluate on 3B. Orange dots denote fitting  
 446 data, and purple crosses indicate the test data. We compare scaling-law predicted loss with actual  
 447 pretraining loss of architectures and we observe that fitting the scaling laws with only 1B model data  
 448 yields lower MSE and higher Spearman correlation for the 3B model loss prediction.  
 449

461 3B across all settings, achieving up to 47% higher throughput with SGLang on H200. This demon-  
 462 strates that the efficiency gains transfer across serving stacks and hardware platforms. Detailed  
 463 throughput statistics are provided in Table 6.

464 **Ablation of fitting data strategy.** While we adopt a progressive strategy for selecting fitting data  
 465 across tasks (§4), results from small models (e.g., 80M) may not reliably predict behaviors at larger  
 466 scales such as 3B. To assess this, we fit the conditional scaling law for the 3B model using only the  
 467 1B variants. As shown in Figure 7, fitting with 1B data yields lower MSE and higher Spearman  
 468 correlation when predicting 3B behavior, suggesting that the law’s coefficients shift with model  
 469 size. We therefore refit the law with multiplicative calibration using only the 1B variants, yielding  
 470 the coefficients  $a_0 = 2.319$ ,  $a_1 = 0.238$ ,  $a_2 = 0.0176$ ,  $b_0 = 0.5104$ ,  $b_1 = 0.0051$ , and  $b_2 = 0.0062$ .  
 471

472 This produces an alternative optimal configuration for the 3B model, with  $d_{\text{model}}/\sqrt{N} = 0.074$  and  
 473  $r = 1.229$ . We train a 3B model (Panda-3B<sup>◦</sup>) under this configuration on 100B tokens and compare  
 474 it with both LLaMA-3.2-3B and Panda-3B (fitted from 80M, 145M, 297M, and 1B data). As shown  
 475 in Table 2, Panda-3B<sup>◦</sup> achieves a lower training loss and comparable downstream accuracy to Panda-  
 476 3B, with detailed results given in Appendix K. These findings suggest that when scaling up, it is  
 477 often sufficient, and sometimes preferable, to fit the law using models within a closer size range to  
 478 the target, such as about one third of its scale.

## 6 RELATED WORK

482 **Large Language Models.** Transformers Vaswani (2017) have shown strong performance across  
 483 diverse downstream tasks, such as text classification Wang (2018); Sarlin et al. (2020), mathematical  
 484 reasoning Cobbe et al. (2021); Hendrycks et al. (2021), and code generation Chen et al. (2021);  
 485 Austin et al. (2021); Jain et al. (2024). The Transformer architecture serves as the foundation for  
 486 many leading large language models, including GPT Brown et al. (2020); Achiam et al. (2023),

486 LLaMA Touvron et al. (2023), Gemma Team et al. (2024a), Qwen Yang et al. (2025), Kimi Team  
 487 et al. (2025), and DeepSeek Liu et al. (2024a); Guo et al. (2025).

489 **Scaling Laws for Language Models.** Scaling laws are powerful tools to predict the performance  
 490 of large language models. Existing scaling laws Hoffmann et al. (2022); Muennighoff et al. (2023);  
 491 Sardana et al. (2023); Kumar et al. (2024); Gadre et al. (2024); Ruan et al. (2024) characterize how  
 492 model performance varies with model size, dataset size, data quality, and compute budget. With the  
 493 rise of Mixture-of-Experts (MoE) Shazeer et al. (2017); Guo et al. (2025), a powerful architecture  
 494 for large language models, recent studies Krajewski et al. (2024); Abnar et al. (2025) extend scaling  
 495 laws to account for the number of experts, expert granularity, active parameters, and sparsity.

496 **Serving Systems.** Due to the increased inference cost, many inference systems have been de-  
 497 veloped to speed up model serving Yu et al. (2022); Kwon et al. (2023); Zheng et al. (2023); Ye  
 498 et al. (2025). Specifically, vLLM Kwon et al. (2023) proposes PagedAttention to manage KV cache  
 499 memory more effectively, thereby improving throughput. Similarly, SGLang Zheng et al. (2023)  
 500 introduces RadixAttention to achieve higher throughput and lower latency.

501 **Inference-Efficient Model Design.** Efforts to improve the inference efficiency of large language  
 502 models generally fall into two categories: one line of work investigates the trade-offs across differ-  
 503 ent model configurations Alabdulmohsin et al. (2023); Bian et al. (2025), while the other focuses  
 504 on designing more efficient model architectures Xiao et al. (2023); Gu & Dao (2023); Gao et al.  
 505 (2024b); Jiang et al. (2024); Liu et al. (2024b); Dao & Gu (2024); Xiao et al. (2024); Yuan et al.  
 506 (2025); Chandrasegaran et al. (2025).

## 510 7 LIMITATIONS AND FUTURE WORK

511 While our team has made notable progress, several open challenges remain that offer promising  
 512 directions for future research. First, due to limitations in resources and time, our evaluation does  
 513 not extend to 7B models. Second, our analysis is restricted to dense models, and it remains unclear  
 514 whether the results extend to Mixture of Experts (MoE) architectures Shazeer et al. (2017). While  
 515 we report inference efficiency measurements for MoE models under varying architectural choices in  
 516 Appendix L, we have not yet established scaling laws for MoE architectures. Finally, our analysis is  
 517 limited to pre-training, and it remains unclear how the results would change under post-training.

## 521 8 CONCLUSION

522 This work explores the trade-off between model accuracy and inference cost under a fixed training  
 523 budget. We begin by demonstrating how architectural choices influence both inference throughput  
 524 and model accuracy. Building on this, we extend Chinchilla scaling laws to incorporate architectural  
 525 factors and propose a two-step conditional framework for optimal architecture search: (i) train small  
 526 models to fit the conditional scaling law (Eq. 3), and (ii) solve Eq. 4 for the predicted optimal  
 527 architecture, followed by a local search over GQA to maximize inference efficiency. Using the fitted  
 528 scaling laws and our framework, we trained models up to 3B parameters, achieving up to 42% higher  
 529 inference throughput and 2.1% accuracy gains across nine downstream tasks. In Table 7 and Table 8  
 530 of Appendix G, we compare design choices across existing open-source models at the 1B and 3B  
 531 scales, further underscoring the need for our inference-efficient accurate model designs.

## 533 REPRODUCIBILITY STATEMENT

535 All experiments in this work were conducted using publicly available frameworks. Section 4  
 536 provides details of our training, inference, and evaluation setups. In particular, we used  
 537 Megatron-LM (Shoeybi et al., 2019) for model training, vLLM (Kwon et al., 2023) for efficient  
 538 inference, and lm-eval-harness (Gao et al., 2024a) for standardized evaluations. To facilitate  
 539 reproducibility, we will release configuration files and scripts.

540 REFERENCES  
541

- 542 Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar,  
543 Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, et al. Phi-4 techni-  
544 cal report. *arXiv preprint arXiv:2412.08905*, 2024.
- 545 Samira Abnar, Harshay Shah, Dan Busbridge, Alaaeldin Mohamed Elnouby Ali, Josh Susskind,  
546 and Vimal Thilak. Parameters vs flops: Scaling laws for optimal sparsity for mixture-of-experts  
547 language models. *arXiv preprint arXiv:2501.12370*, 2025.
- 548 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-  
549 man, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical  
550 report. *arXiv preprint arXiv:2303.08774*, 2023.
- 551 Muhammad Adnan, Akhil Arunkumar, Gaurav Jain, Prashant J Nair, Ilya Soloveychik, and Pu-  
552 rushotham Kamath. Keyformer: Kv cache reduction through key tokens selection for efficient  
553 generative inference. *Proceedings of Machine Learning and Systems*, 6:114–127, 2024.
- 554 Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit  
555 Sanghi. Gqa: Training generalized multi-query transformer models from multi-head check-  
556 points. *arXiv preprint arXiv:2305.13245*, 2023.
- 557 Ibrahim M Alabdulmohsin, Xiaohua Zhai, Alexander Kolesnikov, and Lucas Beyer. Getting vit in  
558 shape: Scaling laws for compute-optimal model design. *Advances in Neural Information Pro-  
559 cessing Systems*, 36:16406–16425, 2023.
- 560 Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan,  
561 Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language  
562 models. *arXiv preprint arXiv:2108.07732*, 2021.
- 563 Song Bian, Minghao Yan, and Shivararam Venkataraman. Scaling inference-efficient language mod-  
564 els. *arXiv preprint arXiv:2501.18107*, 2025.
- 565 Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric  
566 Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al.  
567 Pythia: A suite for analyzing large language models across training and scaling. In *International  
568 Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.
- 569 Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical com-  
570 monsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*,  
571 volume 34, pp. 7432–7439, 2020.
- 572 Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré, and  
573 Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling.  
574 *arXiv preprint arXiv:2407.21787*, 2024.
- 575 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
576 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
577 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 578 Keshigeyan Chandrasegaran, Michael Poli, Daniel Y Fu, Dongjun Kim, Lea M Hadzic, Manling  
579 Li, Agrim Gupta, Stefano Massaroli, Azalia Mirhoseini, Juan Carlos Niebles, et al. Exploring  
580 diffusion transformer designs via grafting. *arXiv preprint arXiv:2506.05340*, 2025.
- 581 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared  
582 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large  
583 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 584 Mengzhao Chen, Chaoyi Zhang, Jing Liu, Yutao Zeng, Zeyue Xue, Zhiheng Liu, Yunshui Li, Jin  
585 Ma, Jie Huang, Xun Zhou, et al. Scaling law for quantization-aware training. *arXiv preprint  
586 arXiv:2505.14302*, 2025.
- 587 Andrew A Chien, Liuzixuan Lin, Hai Nguyen, Varsha Rao, Tristan Sharma, and Rajini Wijayawar-  
588 dana. Reducing the carbon impact of generative ai inference (today and in 2035). In *Proceedings  
589 of the 2nd workshop on sustainable computer systems*, pp. 1–7, 2023.

- 594 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and  
 595 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.  
 596 *arXiv preprint arXiv:1803.05457*, 2018.
- 597
- 598 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
 599 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to  
 600 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 601 Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through  
 602 structured state space duality. *arXiv preprint arXiv:2405.21060*, 2024.
- 603
- 604 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
 605 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
 606 *arXiv preprint arXiv:2407.21783*, 2024.
- 607 Samir Yitzhak Gadre, Georgios Smyrnis, Vaishaal Shankar, Suchin Gururangan, Mitchell Worts-  
 608 man, Rulin Shao, Jean Mercat, Alex Fang, Jeffrey Li, Sedrick Keh, et al. Language models scale  
 609 reliably with over-training and on downstream tasks. *arXiv preprint arXiv:2403.08540*, 2024.
- 610
- 611 Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Fos-  
 612 ter, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muen-  
 613 nighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang  
 614 Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. The language model  
 615 evaluation harness, 07 2024a. URL <https://zenodo.org/records/12608602>.
- 616
- 617 Yizhao Gao, Zhichen Zeng, Dayou Du, Shijie Cao, Peiyuan Zhou, Jiaxing Qi, Junjie Lai, Hayden  
 618 Kwok-Hay So, Ting Cao, Fan Yang, et al. Seerattention: Learning intrinsic sparse attention in  
 619 your llms. *arXiv preprint arXiv:2410.13276*, 2024b.
- 620
- 621 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*  
 622 *preprint arXiv:2312.00752*, 2023.
- 623
- 624 Xinyu Guan, Li Lyra Zhang, Yifei Liu, Ning Shang, Youran Sun, Yi Zhu, Fan Yang, and Mao Yang.  
 625 rstar-math: Small llms can master math reasoning with self-evolved deep thinking. *arXiv preprint*  
 626 *arXiv:2501.04519*, 2025.
- 627
- 628 Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,  
 629 Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms  
 630 via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- 631
- 632 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,  
 633 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *arXiv*  
 634 *preprint arXiv:2103.03874*, 2021.
- 635
- 636 Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza  
 637 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Train-  
 638 ing compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- 639
- 640 Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando  
 641 Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free  
 642 evaluation of large language models for code. *arXiv preprint arXiv:2403.07974*, 2024.
- 643
- 644 Huiqiang Jiang, Yucheng Li, Chengruidong Zhang, Qianhui Wu, Xufang Luo, Surin Ahn, Zhenhua  
 645 Han, Amir H Abdi, Dongsheng Li, Chin-Yew Lin, et al. Minference 1.0: Accelerating pre-filling  
 646 for long-context llms via dynamic sparse attention. *arXiv preprint arXiv:2407.02490*, 2024.
- 647
- 648 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child,  
 649 Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language  
 650 models. *arXiv preprint arXiv:2001.08361*, 2020.
- 651
- 652 Jakub Krajewski, Jan Ludziejewski, Kamil Adamczewski, Maciej Pióro, Michał Krutul, Szymon  
 653 Antoniak, Kamil Ciebiera, Krystian Król, Tomasz Odrzygóźdż, Piotr Sankowski, et al. Scaling  
 654 laws for fine-grained mixture of experts. *arXiv preprint arXiv:2402.07871*, 2024.

- 648 Tanishq Kumar, Zachary Ankner, Benjamin F Spector, Blake Bordelon, Niklas Muennighoff, Man-  
 649 sheej Paul, Cengiz Pehlevan, Christopher Ré, and Aditi Raghunathan. Scaling laws for precision.  
 650 *arXiv preprint arXiv:2411.04330*, 2024.
- 651
- 652 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph  
 653 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model  
 654 serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Prin-*  
 655 *ciples*, pp. 611–626, 2023.
- 656 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,  
 657 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint*  
 658 *arXiv:2412.19437*, 2024a.
- 659
- 660 Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao,  
 661 Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint*  
 662 *arXiv:2412.19437*, 2024b.
- 663 Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat Phatale, Harsh Lara, Yunxuan Li, Lei Shu, Yun  
 664 Zhu, Lei Meng, Jiao Sun, et al. Improve mathematical reasoning in language models by automated  
 665 process supervision. *arXiv preprint arXiv:2406.06592*, 2024.
- 666
- 667 Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct  
 668 electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*,  
 669 2018.
- 670
- 671 Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra  
 672 Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. Scaling data-constrained language  
 673 models. *Advances in Neural Information Processing Systems*, 36:50358–50376, 2023.
- 674
- 675 Aashiq Muhamed, Christian Bock, Rahul Solanki, Youngsuk Park, Yida Wang, and Jun Huan. Train-  
 676 ing large-scale foundation models on emerging ai chips. In *Proceedings of the 29th ACM SIGKDD*  
 677 *Conference on Knowledge Discovery and Data Mining*, pp. 5821–5822, 2023.
- 678
- 679 Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi,  
 680 Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The lambada dataset:  
 681 Word prediction requiring a broad discourse context. *arXiv preprint arXiv:1606.06031*, 2016.
- 682
- 683 Youngsuk Park, Kailash Budhathoki, Liangfu Chen, Jonas M Kübler, Jiaji Huang, Matthäus Klein-  
 684 dessner, Jun Huan, Volkan Cevher, Yida Wang, and George Karypis. Inference optimization of  
 685 foundation models on ai accelerators. In *Proceedings of the 30th ACM SIGKDD Conference on*  
 686 *Knowledge Discovery and Data Mining*, pp. 6605–6615, 2024.
- 687
- 688 Jackson Petty, Sjoerd van Steenkiste, Ishita Dasgupta, Fei Sha, Dan Garrette, and Tal Linzen. The  
 689 impact of depth on compositional generalization in transformer language models. *arXiv preprint*  
 690 *arXiv:2310.19956*, 2023.
- 691
- 692 Zhenting Qi, Mingyuan Ma, Jiahang Xu, Li Lyra Zhang, Fan Yang, and Mao Yang. Mutual reason-  
 693 ing makes smaller llms stronger problem-solvers. *arXiv preprint arXiv:2408.06195*, 2024.
- 694
- 695 Siva Reddy, Danqi Chen, and Christopher D Manning. Coqa: A conversational question answering  
 696 challenge. *Transactions of the Association for Computational Linguistics*, 7:249–266, 2019.
- 697
- 698 Yangjun Ruan, Chris J Maddison, and Tatsunori Hashimoto. Observational scaling laws and the  
 699 predictability of language model performance. *arXiv preprint arXiv:2405.10938*, 2024.
- 700
- 701 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adver-  
 702 sarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- 703
- Nikhil Sardana, Jacob Portes, Sasha Doubov, and Jonathan Frankle. Beyond chinchilla-optimal:  
 Accounting for inference in language model scaling laws. *arXiv preprint arXiv:2401.00448*,  
 2023.

- 702 Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue:  
 703 Learning feature matching with graph neural networks. In *Proceedings of the IEEE/CVF confer-  
 704 ence on computer vision and pattern recognition*, pp. 4938–4947, 2020.
- 705 Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton,  
 706 and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer.  
 707 *arXiv preprint arXiv:1701.06538*, 2017.
- 708 Mohammad Shoeybi, Mostafa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan  
 709 Catanzaro. Megatron-LM: Training multi-billion parameter language models using model par-  
 710 allelism. *arXiv preprint arXiv:1909.08053*, 2019.
- 711 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling LLM test-time compute optimally  
 712 can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- 713 Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur,  
 714 Ben Beglin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, et al. Dolma: An open corpus of  
 715 three trillion tokens for language model pretraining research. *arXiv preprint arXiv:2402.00159*,  
 716 2024.
- 717 Charles Spearman. The proof and measurement of association between two things. 1961.
- 718 Yi Tay, Mostafa Dehghani, Jinfeng Rao, William Fedus, Samira Abnar, Hyung Won Chung, Sharan  
 719 Narang, Dani Yogatama, Ashish Vaswani, and Donald Metzler. Scale efficiently: Insights from  
 720 pre-training and fine-tuning transformers. *arXiv preprint arXiv:2109.10686*, 2021.
- 721 Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya  
 722 Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open  
 723 models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024a.
- 724 Gemma Team, Morgane Rivière, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhu-  
 725 patiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma  
 726 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024b.
- 727 Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen,  
 728 Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. *arXiv  
 729 preprint arXiv:2507.20534*, 2025.
- 730 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
 731 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and  
 732 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 733 A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- 734 Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head  
 735 self-attention: Specialized heads do the heavy lifting, the rest can be pruned. *arXiv preprint  
 736 arXiv:1905.09418*, 2019.
- 737 Alex Wang. Glue: A multi-task benchmark and analysis platform for natural language understand-  
 738 ing. *arXiv preprint arXiv:1804.07461*, 2018.
- 739 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yo-  
 740 gatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language  
 741 models. *arXiv preprint arXiv:2206.07682*, 2022.
- 742 Johannes Welbl, Nelson F Liu, and Matt Gardner. Crowdsourcing multiple choice science questions.  
 743 *arXiv preprint arXiv:1707.06209*, 2017.
- 744 Carole-Jean Wu, Bilge Acun, Ramya Raghavendra, and Kim Hazelwood. Beyond efficiency: Scal-  
 745 ing ai sustainably. *IEEE Micro*, 44(5):37–46, 2024.
- 746 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming  
 747 language models with attention sinks. *arXiv preprint arXiv:2309.17453*, 2023.

- 756   Guangxuan Xiao, Jiaming Tang, Jingwei Zuo, Junxian Guo, Shang Yang, Haotian Tang, Yao Fu,  
 757   and Song Han. Duoattention: Efficient long-context llm inference with retrieval and streaming  
 758   heads. *arXiv preprint arXiv:2410.10819*, 2024.
- 759   An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,  
 760   Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint*  
 761   *arXiv:2505.09388*, 2025.
- 762   Zihao Ye, Lequn Chen, Ruihang Lai, Wuwei Lin, Yineng Zhang, Stephanie Wang, Tianqi Chen,  
 763   Baris Kasikci, Vinod Grover, Arvind Krishnamurthy, et al. Flashinfer: Efficient and customizable  
 764   attention engine for llm inference serving. *arXiv preprint arXiv:2501.01005*, 2025.
- 765   Gyeong-In Yu, Joo Seong Jeong, Geon-Woo Kim, Soojeong Kim, and Byung-Gon Chun. Orca: A  
 766   distributed serving system for {Transformer-Based} generative models. In *16th USENIX Symposium on*  
 767   *Operating Systems Design and Implementation (OSDI 22)*, pp. 521–538, 2022.
- 768   Jingyang Yuan, Huazuo Gao, Damai Dai, Junyu Luo, Liang Zhao, Zhengyan Zhang, Zhenda Xie,  
 769   YX Wei, Lean Wang, Zhiping Xiao, et al. Native sparse attention: Hardware-aligned and natively  
 770   trainable sparse attention. *arXiv preprint arXiv:2502.11089*, 2025.
- 771   Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a ma-  
 772   chine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.
- 773   Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small  
 774   language model. *arXiv preprint arXiv:2401.02385*, 2024.
- 775   Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Jeff Huang, Chuyue Sun, Cody\_Hao Yu, Shiyi Cao,  
 776   Christos Kozyrakis, Ion Stoica, Joseph E Gonzalez, et al. Efficiently programming large language  
 777   models using sclang. 2023.
- 778  
 779  
 780  
 781  
 782  
 783  
 784  
 785  
 786  
 787  
 788  
 789  
 790  
 791  
 792  
 793  
 794  
 795  
 796  
 797  
 798  
 799  
 800  
 801  
 802  
 803  
 804  
 805  
 806  
 807  
 808  
 809

## A LLM USAGE

We used an LLM to improve the writing by correcting grammar in our draft. It was not used to generate research ideas.

## B OPEN-WEIGHTED MODEL ARCHITECTURES

Table 3 presents an overview of the open-weight model architectures utilized in this paper.

**Table 3: Open-Weighted Model Architectures:** We list the architectural configurations of all models used in this paper.  $n_{\text{layers}}$  is the number of layers,  $d_{\text{model}}$  is the hidden size,  $n_{\text{heads}}$  is the number of attention heads, and  $f_{\text{size}}$  is the intermediate size.

| Model Name   | $n_{\text{layers}}$ | $d_{\text{model}}$ | $n_{\text{heads}}$ | $f_{\text{size}}$ | GQA |
|--------------|---------------------|--------------------|--------------------|-------------------|-----|
| Qwen2.5-1.5B | 28                  | 1536               | 12                 | 8960              | 6   |
| Qwen3-0.6B   | 28                  | 1024               | 16                 | 3072              | 2   |

## C MODEL ARCHITECTURES

Table 4 provides an overview of the model architectures, all configured with  $GQA = 4$  and employing LLaMA-3.2 as the tokenizer.

Table 4: **Model Architectures:** We list the architectural configurations of all models trained in this paper.  $N_{\text{non-embed}}$  is the total number of non-embedding parameters,  $n_{\text{layers}}$  is the number of layers,  $d_{\text{model}}$  is the hidden size,  $n_{\text{heads}}$  is the number of attention heads,  $f_{\text{size}}$  is the intermediate size, and  $r_{\text{mlp}/\text{attn}}$  is the MLP-to-attention ratio.

| $N_{\text{non-embed}}$ | Variant | $n_{\text{layers}}$ | $d_{\text{model}}$ | $n_{\text{heads}}$ | $f_{\text{size}}$ | $d_{\text{model}}/\sqrt{N}$ | $r_{\text{mlp/attn}}$ |
|------------------------|---------|---------------------|--------------------|--------------------|-------------------|-----------------------------|-----------------------|
| 80M                    | v1      | 12                  | 768                | 16                 | 2048              | 0.086                       | 2.40                  |
| 80M                    | v2      | 12                  | 768                | 4                  | 2688              | 0.086                       | 12.6                  |
| 80M                    | v3      | 12                  | 768                | 8                  | 2560              | 0.085                       | 6.00                  |
| 80M                    | v4      | 12                  | 768                | 24                 | 1536              | 0.087                       | 1.20                  |
| 80M                    | v5      | 12                  | 768                | 32                 | 1152              | 0.086                       | 0.68                  |
| 80M                    | v6      | 12                  | 768                | 40                 | 768               | 0.086                       | 0.36                  |
| 80M                    | v7      | 12                  | 768                | 48                 | 256               | 0.087                       | 0.10                  |
| 80M                    | v8      | 12                  | 384                | 32                 | 4096              | 0.043                       | 2.40                  |
| 80M                    | v9      | 12                  | 384                | 8                  | 5376              | 0.043                       | 12.6                  |
| 80M                    | v10     | 12                  | 384                | 16                 | 5120              | 0.042                       | 6.00                  |
| 80M                    | v11     | 12                  | 384                | 48                 | 3072              | 0.044                       | 1.20                  |
| 80M                    | v12     | 12                  | 384                | 64                 | 2304              | 0.043                       | 0.68                  |
| 80M                    | v13     | 12                  | 384                | 80                 | 1536              | 0.043                       | 0.36                  |
| 80M                    | v14     | 12                  | 384                | 96                 | 512               | 0.044                       | 0.10                  |
| 80M                    | v15     | 12                  | 1536               | 8                  | 1024              | 0.171                       | 2.40                  |
| 80M                    | v16     | 12                  | 1536               | 4                  | 1280              | 0.169                       | 6.00                  |
| 80M                    | v17     | 12                  | 1536               | 12                 | 768               | 0.174                       | 1.20                  |
| 80M                    | v18     | 12                  | 1536               | 16                 | 640               | 0.169                       | 0.75                  |
| 80M                    | v19     | 12                  | 1536               | 20                 | 384               | 0.171                       | 0.36                  |
| 80M                    | v20     | 12                  | 1536               | 24                 | 128               | 0.174                       | 0.10                  |
| 80M                    | v21     | 12                  | 512                | 24                 | 3072              | 0.057                       | 2.40                  |
| 80M                    | v22     | 12                  | 512                | 12                 | 3840              | 0.056                       | 6.00                  |
| 80M                    | v23     | 12                  | 512                | 16                 | 3584              | 0.057                       | 4.20                  |
| 80M                    | v24     | 12                  | 512                | 36                 | 2304              | 0.058                       | 1.20                  |
| 80M                    | v25     | 12                  | 512                | 48                 | 1792              | 0.057                       | 0.70                  |
| 80M                    | v26     | 12                  | 512                | 60                 | 1152              | 0.057                       | 0.36                  |

| 864 | $N_{\text{non-embed}}$ | Variant | $n_{\text{layers}}$ | $d_{\text{model}}$ | $n_{\text{heads}}$ | $f_{\text{size}}$ | $d_{\text{model}}/\sqrt{N}$ | $r_{\text{mlp}/\text{attn}}$ |
|-----|------------------------|---------|---------------------|--------------------|--------------------|-------------------|-----------------------------|------------------------------|
| 865 | 80M                    | v27     | 12                  | 512                | 72                 | 384               | 0.058                       | 0.10                         |
| 866 | 80M                    | v28     | 12                  | 1024               | 12                 | 1536              | 0.114                       | 2.40                         |
| 867 | 80M                    | v29     | 12                  | 1024               | 8                  | 1792              | 0.113                       | 4.20                         |
| 868 | 80M                    | v30     | 12                  | 1024               | 16                 | 1280              | 0.115                       | 1.50                         |
| 869 | 80M                    | v31     | 12                  | 1024               | 24                 | 896               | 0.114                       | 0.70                         |
| 870 | 80M                    | v32     | 12                  | 1024               | 36                 | 256               | 0.114                       | 0.13                         |
| 871 | 80M                    | v33     | 12                  | 2048               | 4                  | 896               | 0.226                       | 4.20                         |
| 872 | 80M                    | v34     | 12                  | 2048               | 8                  | 640               | 0.231                       | 1.50                         |
| 873 | 80M                    | v35     | 12                  | 2048               | 16                 | 256               | 0.226                       | 0.30                         |
| 874 | 80M                    | v48     | 12                  | 768                | 20                 | 1792              | 0.086                       | 1.68                         |
| 875 | 80M                    | v49     | 12                  | 768                | 28                 | 1408              | 0.086                       | 0.94                         |
| 876 | 80M                    | v50     | 12                  | 384                | 40                 | 3584              | 0.043                       | 1.68                         |
| 877 | 80M                    | v51     | 12                  | 384                | 52                 | 3072              | 0.043                       | 1.11                         |
| 878 | 80M                    | v52     | 12                  | 384                | 56                 | 2816              | 0.043                       | 0.94                         |
| 879 | 80M                    | v53     | 12                  | 384                | 60                 | 2560              | 0.043                       | 0.80                         |
| 880 | 80M                    | v54     | 12                  | 512                | 32                 | 2560              | 0.058                       | 1.50                         |
| 881 | 80M                    | v55     | 12                  | 512                | 40                 | 2176              | 0.057                       | 1.02                         |
| 882 | 80M                    | v56     | 12                  | 512                | 44                 | 1920              | 0.058                       | 0.82                         |
| 883 | 80M                    | v57     | 12                  | 1024               | 20                 | 1152              | 0.113                       | 1.08                         |
| 884 | 145M                   | v1      | 12                  | 1024               | 16                 | 3072              | 0.085                       | 3.60                         |
| 885 | 145M                   | v2      | 12                  | 1024               | 8                  | 3584              | 0.084                       | 8.40                         |
| 886 | 145M                   | v3      | 12                  | 1024               | 24                 | 2560              | 0.086                       | 2.00                         |
| 887 | 145M                   | v4      | 12                  | 1024               | 32                 | 2304              | 0.084                       | 1.35                         |
| 888 | 145M                   | v5      | 12                  | 1024               | 40                 | 1792              | 0.085                       | 0.84                         |
| 889 | 145M                   | v6      | 12                  | 1024               | 48                 | 1280              | 0.086                       | 0.50                         |
| 890 | 145M                   | v7      | 12                  | 1024               | 64                 | 512               | 0.085                       | 0.15                         |
| 891 | 145M                   | v8      | 12                  | 512                | 32                 | 6144              | 0.043                       | 3.60                         |
| 892 | 145M                   | v9      | 12                  | 512                | 16                 | 7168              | 0.042                       | 8.40                         |
| 893 | 145M                   | v10     | 12                  | 512                | 48                 | 5120              | 0.043                       | 2.00                         |
| 894 | 145M                   | v11     | 12                  | 512                | 64                 | 4608              | 0.042                       | 1.35                         |
| 895 | 145M                   | v12     | 12                  | 512                | 80                 | 3584              | 0.043                       | 0.84                         |
| 896 | 145M                   | v13     | 12                  | 512                | 96                 | 2560              | 0.043                       | 0.50                         |
| 897 | 145M                   | v14     | 12                  | 512                | 128                | 1024              | 0.043                       | 0.15                         |
| 898 | 145M                   | v15     | 12                  | 2048               | 8                  | 1536              | 0.170                       | 3.60                         |
| 899 | 145M                   | v16     | 12                  | 2048               | 4                  | 1792              | 0.168                       | 8.40                         |
| 900 | 145M                   | v17     | 12                  | 2048               | 12                 | 1280              | 0.172                       | 2.00                         |
| 901 | 145M                   | v18     | 12                  | 2048               | 16                 | 1152              | 0.168                       | 1.35                         |
| 902 | 145M                   | v19     | 12                  | 2048               | 20                 | 896               | 0.170                       | 0.84                         |
| 903 | 145M                   | v20     | 12                  | 2048               | 24                 | 640               | 0.172                       | 0.50                         |
| 904 | 145M                   | v21     | 12                  | 2048               | 32                 | 256               | 0.170                       | 0.15                         |
| 905 | 145M                   | v22     | 12                  | 768                | 24                 | 3840              | 0.065                       | 3.00                         |
| 906 | 145M                   | v23     | 12                  | 768                | 32                 | 3584              | 0.063                       | 2.10                         |
| 907 | 145M                   | v24     | 12                  | 768                | 40                 | 3072              | 0.064                       | 1.44                         |
| 908 | 145M                   | v25     | 12                  | 768                | 48                 | 2560              | 0.065                       | 1.00                         |
| 909 | 145M                   | v26     | 12                  | 768                | 56                 | 2304              | 0.063                       | 0.77                         |
| 910 | 145M                   | v27     | 12                  | 768                | 64                 | 1792              | 0.064                       | 0.53                         |
| 911 | 145M                   | v28     | 12                  | 1536               | 12                 | 1920              | 0.129                       | 3.00                         |
| 912 | 145M                   | v29     | 12                  | 1536               | 16                 | 1792              | 0.127                       | 2.10                         |
| 913 | 145M                   | v30     | 12                  | 1536               | 20                 | 1536              | 0.128                       | 1.44                         |
| 914 | 145M                   | v31     | 12                  | 1536               | 24                 | 1280              | 0.129                       | 1.00                         |
| 915 | 145M                   | v32     | 12                  | 1536               | 28                 | 1152              | 0.127                       | 0.77                         |
| 916 | 145M                   | v33     | 12                  | 1536               | 32                 | 896               | 0.128                       | 0.53                         |
| 917 | 145M                   | v34     | 12                  | 4096               | 4                  | 768               | 0.340                       | 3.60                         |
| 918 | 145M                   | v35     | 12                  | 4096               | 16                 | 128               | 0.340                       | 0.15                         |
| 919 | 145M                   | v48     | 12                  | 1024               | 28                 | 2368              | 0.086                       | 1.59                         |
| 920 | 145M                   | v49     | 12                  | 1024               | 36                 | 2048              | 0.085                       | 1.07                         |
| 921 | 145M                   | v50     | 12                  | 512                | 52                 | 5120              | 0.042                       | 1.85                         |

| 918 | $N_{\text{non-embed}}$ | Variant | $n_{\text{layers}}$ | $d_{\text{model}}$ | $n_{\text{heads}}$ | $f_{\text{size}}$ | $d_{\text{model}}/\sqrt{N}$ | $r_{\text{mlp/attn}}$ |
|-----|------------------------|---------|---------------------|--------------------|--------------------|-------------------|-----------------------------|-----------------------|
| 919 | 145M                   | v51     | 12                  | 512                | 60                 | 4800              | 0.042                       | 1.50                  |
| 920 | 145M                   | v52     | 12                  | 512                | 68                 | 4224              | 0.043                       | 1.16                  |
| 921 | 145M                   | v53     | 12                  | 512                | 72                 | 3968              | 0.043                       | 1.03                  |
| 922 | 145M                   | v54     | 12                  | 768                | 44                 | 2944              | 0.063                       | 1.25                  |
| 923 | 145M                   | v55     | 12                  | 768                | 52                 | 2432              | 0.064                       | 0.88                  |
| 924 | 297M                   | v1      | 12                  | 1536               | 24                 | 4096              | 0.089                       | 3.20                  |
| 925 | 297M                   | v2      | 12                  | 1536               | 8                  | 4864              | 0.090                       | 11.4                  |
| 926 | 297M                   | v3      | 12                  | 1536               | 16                 | 4608              | 0.088                       | 5.40                  |
| 927 | 297M                   | v4      | 12                  | 1536               | 32                 | 3584              | 0.090                       | 2.10                  |
| 928 | 297M                   | v5      | 12                  | 1536               | 48                 | 2816              | 0.089                       | 1.10                  |
| 929 | 297M                   | v6      | 12                  | 1536               | 64                 | 2048              | 0.088                       | 0.60                  |
| 930 | 297M                   | v7      | 12                  | 1536               | 80                 | 1024              | 0.090                       | 0.24                  |
| 931 | 297M                   | v8      | 12                  | 768                | 48                 | 8192              | 0.045                       | 3.20                  |
| 932 | 297M                   | v9      | 12                  | 768                | 16                 | 9728              | 0.045                       | 11.4                  |
| 933 | 297M                   | v10     | 12                  | 768                | 32                 | 9216              | 0.044                       | 5.40                  |
| 934 | 297M                   | v11     | 12                  | 768                | 64                 | 7168              | 0.045                       | 2.10                  |
| 935 | 297M                   | v12     | 12                  | 768                | 96                 | 5632              | 0.045                       | 1.10                  |
| 936 | 297M                   | v13     | 12                  | 768                | 128                | 4096              | 0.044                       | 0.60                  |
| 937 | 297M                   | v14     | 12                  | 768                | 160                | 2048              | 0.045                       | 0.24                  |
| 938 | 297M                   | v15     | 12                  | 3072               | 12                 | 2048              | 0.178                       | 3.20                  |
| 939 | 297M                   | v16     | 12                  | 3072               | 4                  | 2432              | 0.180                       | 11.4                  |
| 940 | 297M                   | v17     | 12                  | 3072               | 8                  | 2304              | 0.177                       | 5.40                  |
| 941 | 297M                   | v18     | 12                  | 3072               | 16                 | 1792              | 0.180                       | 2.10                  |
| 942 | 297M                   | v19     | 12                  | 3072               | 24                 | 1408              | 0.178                       | 1.10                  |
| 943 | 297M                   | v20     | 12                  | 3072               | 32                 | 1024              | 0.177                       | 0.60                  |
| 944 | 297M                   | v21     | 12                  | 3072               | 40                 | 512               | 0.180                       | 0.24                  |
| 945 | 297M                   | v22     | 12                  | 1024               | 36                 | 6144              | 0.059                       | 3.20                  |
| 946 | 297M                   | v23     | 12                  | 1024               | 12                 | 7296              | 0.060                       | 11.4                  |
| 947 | 297M                   | v24     | 12                  | 1024               | 24                 | 6912              | 0.059                       | 5.40                  |
| 948 | 297M                   | v25     | 12                  | 1024               | 48                 | 5376              | 0.060                       | 2.10                  |
| 949 | 297M                   | v26     | 12                  | 1024               | 72                 | 4224              | 0.059                       | 1.10                  |
| 950 | 297M                   | v27     | 12                  | 1024               | 96                 | 3072              | 0.059                       | 0.60                  |
| 951 | 297M                   | v28     | 12                  | 1024               | 120                | 1536              | 0.060                       | 0.24                  |
| 952 | 297M                   | v29     | 12                  | 2048               | 12                 | 3456              | 0.118                       | 5.40                  |
| 953 | 297M                   | v30     | 12                  | 2048               | 24                 | 2688              | 0.120                       | 2.10                  |
| 954 | 297M                   | v31     | 12                  | 2048               | 48                 | 1536              | 0.118                       | 0.60                  |
| 955 | 297M                   | v32     | 12                  | 2048               | 60                 | 768               | 0.120                       | 0.24                  |
| 956 | 297M                   | v45     | 12                  | 1536               | 40                 | 3200              | 0.089                       | 1.50                  |
| 957 | 297M                   | v46     | 12                  | 1536               | 44                 | 3072              | 0.089                       | 1.31                  |
| 958 | 297M                   | v47     | 12                  | 1536               | 52                 | 2688              | 0.088                       | 0.97                  |
| 959 | 297M                   | v48     | 12                  | 1536               | 56                 | 2432              | 0.089                       | 0.81                  |
| 960 | 297M                   | v49     | 12                  | 768                | 80                 | 6400              | 0.045                       | 1.50                  |
| 961 | 297M                   | v50     | 12                  | 768                | 88                 | 6016              | 0.045                       | 1.28                  |
| 962 | 297M                   | v51     | 12                  | 768                | 104                | 5376              | 0.044                       | 0.97                  |
| 963 | 297M                   | v52     | 12                  | 768                | 112                | 4736              | 0.045                       | 0.79                  |
| 964 | 297M                   | v53     | 12                  | 3072               | 20                 | 1664              | 0.177                       | 1.56                  |
| 965 | 297M                   | v54     | 12                  | 3072               | 28                 | 1152              | 0.180                       | 0.77                  |
| 966 | 297M                   | v55     | 12                  | 1024               | 56                 | 4864              | 0.060                       | 1.63                  |
| 967 | 297M                   | v56     | 12                  | 1024               | 64                 | 4608              | 0.060                       | 1.35                  |
| 968 | 297M                   | v57     | 12                  | 1024               | 80                 | 3840              | 0.059                       | 0.90                  |
| 969 | 297M                   | v58     | 12                  | 1024               | 88                 | 3328              | 0.060                       | 0.71                  |
| 970 | 297M                   | v59     | 12                  | 2048               | 32                 | 2432              | 0.117                       | 1.43                  |
| 971 | 297M                   | v60     | 12                  | 2048               | 36                 | 2048              | 0.120                       | 1.07                  |
| 971 | 1B                     | v1      | 16                  | 2048               | 32                 | 8192              | 0.066                       | 4.80                  |
| 971 | 1B                     | v2      | 16                  | 2048               | 72                 | 5760              | 0.067                       | 1.50                  |

---

| 972 | $N_{\text{non-embed}}$ | Variant | $n_{\text{layers}}$ | $d_{\text{model}}$ | $n_{\text{heads}}$ | $f_{\text{size}}$ | $d_{\text{model}}/\sqrt{N}$ | $r_{\text{mlp/attn}}$ |
|-----|------------------------|---------|---------------------|--------------------|--------------------|-------------------|-----------------------------|-----------------------|
| 973 | 1B                     | v3      | 16                  | 2816               | 92                 | 2432              | 0.089                       | 0.50                  |
| 974 | 1B                     | v4      | 16                  | 2816               | 76                 | 3072              | 0.091                       | 0.76                  |
| 975 | 1B                     | v5      | 16                  | 2816               | 68                 | 3584              | 0.090                       | 0.99                  |
| 976 | 1B                     | v6      | 16                  | 2816               | 60                 | 4096              | 0.090                       | 1.28                  |
| 977 | 1B                     | v7      | 16                  | 2816               | 56                 | 4480              | 0.089                       | 1.50                  |
| 978 | 1B                     | v8      | 16                  | 2816               | 24                 | 6144              | 0.089                       | 4.80                  |
| 979 | 1B                     | v9      | 16                  | 2816               | 48                 | 4736              | 0.090                       | 1.85                  |
| 980 | 1B                     | v10     | 16                  | 2816               | 40                 | 5120              | 0.090                       | 2.40                  |
| 981 | 1B                     | v11     | 16                  | 2816               | 36                 | 5376              | 0.090                       | 2.80                  |
| 982 | 1B                     | v12     | 16                  | 2560               | 64                 | 4480              | 0.082                       | 1.31                  |
| 983 | 1B                     | v13     | 16                  | 2560               | 72                 | 4096              | 0.082                       | 1.07                  |
| 984 | 1B                     | v14     | 16                  | 2560               | 80                 | 3648              | 0.082                       | 0.86                  |
| 985 | 1B                     | v15     | 16                  | 2560               | 56                 | 4864              | 0.082                       | 1.63                  |
| 986 | 1B                     | v16     | 16                  | 2560               | 88                 | 3200              | 0.082                       | 0.68                  |
| 987 | 1B                     | v17     | 16                  | 2560               | 48                 | 5376              | 0.082                       | 2.10                  |

---

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

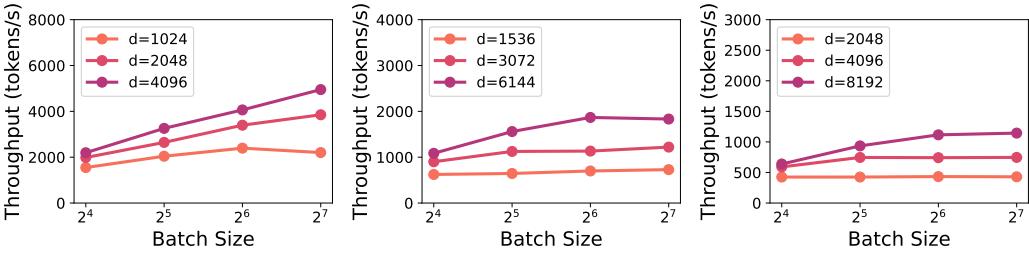
1026 **D HYPER-PARAMETERS**  
10271028 Table 5 lists the detailed hyper-parameters used for training in this paper.  
10291030 **Table 5: Hyper-parameters:** We show the hyper-parameters used for training in this paper.  
1031

|                 | Model Size                                | 80M          | 145M   | 297M   | 1B     | 3B     |
|-----------------|-------------------------------------------|--------------|--------|--------|--------|--------|
| Batch Size      | 256                                       | 256          | 512    | 512    | 512    | 512    |
| Max LR          | 1.5e-3                                    | 1.0e-3       | 8.0e-4 | 6.0e-4 | 6.0e-4 | 6.0e-4 |
| Min LR          |                                           | 0.1 × Max LR |        |        |        |        |
| Optimizer       | AdamW ( $\beta_1 = 0.9, \beta_2 = 0.95$ ) |              |        |        |        |        |
| Weight Decay    |                                           | 0.1          |        |        |        |        |
| Clip Grad Norm  |                                           |              | 1.0    |        |        |        |
| LR Schedule     |                                           |              | Cosine |        |        |        |
| Warmup Steps    |                                           |              | 500    |        |        |        |
| Sequence Length |                                           |              | 2048   |        |        |        |

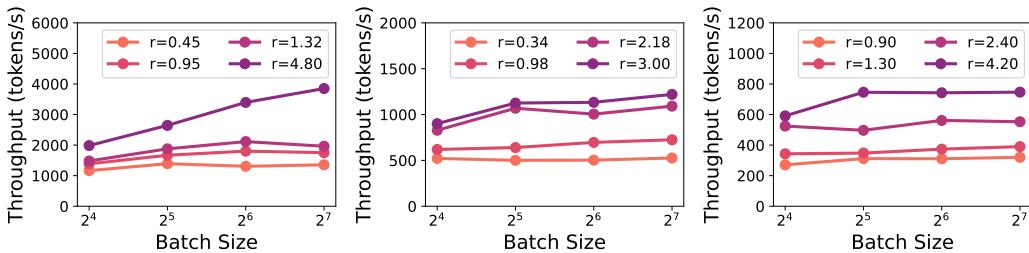
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

## 1080 E ADDITIONAL INFERENCE EVALUATION RESULTS OVER A100 GPUs

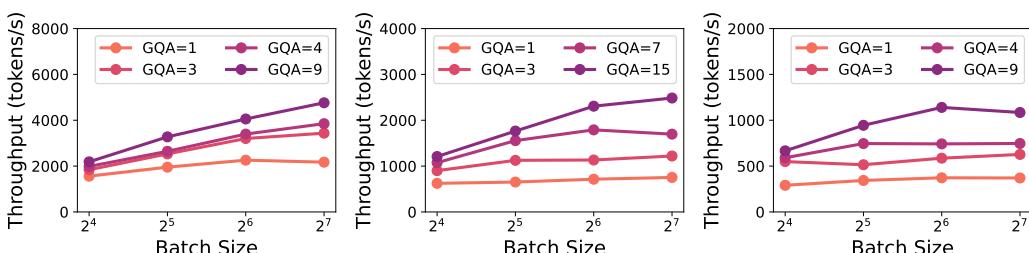
1082 In this section, we present additional inference efficiency results on NVIDIA A100 GPUs. Figure 10  
 1083 presents that, when parameter count, MLP-to-Attention ratio, and hidden size are fixed, increasing  
 1084 GQA leads to higher inference throughput, consistent with the findings of Ainslie et al. (2023). We  
 1085 alter model configurations of LLaMA-3.2-1B, 3B, and LLaMA-3.1-8B in Figure 10. Moreover, we  
 1086 use the SGLang framework Zheng et al. (2023) to benchmark the inference throughput of LLaMA-  
 1087 3.2-1B, LLaMA-3.2-3B, Surefire-1B, and Surefire-3B on a single A100 GPU in Figure 14.



1088 Figure 8: **Hidden size on Inference Throughput:** (left) 1B model variants; (center) 3B model  
 1089 variants; (right) 8B model variants. Across varying batch sizes and model scales, larger hidden  
 1090 sizes yield higher inference throughput under a fixed parameter budget. The legend indicates the hidden  
 1091 size of the models, where  $d = d_{\text{model}}$ .



1102 Figure 9: **MLP-to-Attention ratio on Inference Throughput:** (left) 1B model variants; (center)  
 1103 3B model variants; (right) 8B model variants. Across varying batch sizes and model scales, a larger  
 1104 MLP-to-Attention ratio increases inference throughput under a fixed parameter budget. The legend  
 1105 indicates the MLP-to-Attention ratio of the models, where  $r = r_{\text{mlp}/\text{attn}}$ .



1118 Figure 10: **GQA on Inference Throughput:** (left) 1B model variants; (center) 3B model variants;  
 1119 (right) 8B model variants. This figure shows the impact of GQA on inference throughput. With the  
 1120 total parameter count fixed, hidden size is set to 2048 (1B), 3072 (3B), and 4096 (8B), and the MLP-  
 1121 to-Attention ratio is 4.0, 2.67, and 4.2, respectively. Across varying batch sizes, models with larger  
 1122 GQA achieve higher throughput. All evaluations are performed using the vLLM framework Kwon  
 1123 et al. (2023) on a single NVIDIA Ampere 40GB A100 GPU with 4096 input and 1024 output tokens.

Furthermore, we derive architectural variants by altering the configurations of Qwen3-0.6B, 1.7B, and 4B to investigate the impact of model architectural factors on inference efficiency. The results are shown in Figure 11-13.

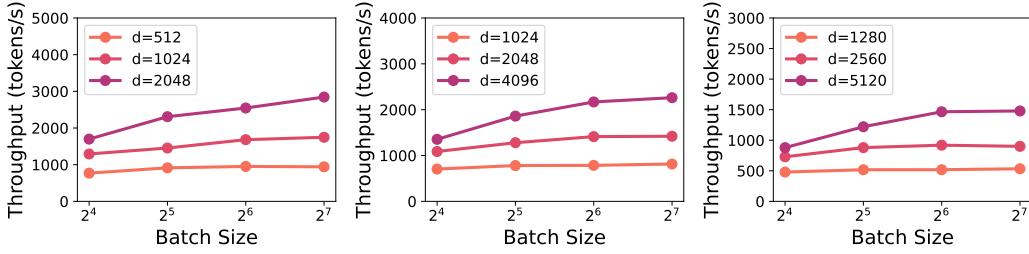


Figure 11: **Hidden size on Inference Throughput (Qwen3):** (left) Qwen3-0.6B model variants; (center) Qwen3-1.7B model variants; (right) Qwen3-4B model variants. Across varying batch sizes and model scales, larger hidden sizes yield higher inference throughput under a fixed parameter budget. The legend indicates the hidden size of the models, where  $d = d_{\text{model}}$ . All evaluations are performed using the vLLM framework Kwon et al. (2023) on a single NVIDIA Ampere 40GB A100 GPU with 4096 input and 1024 output tokens.

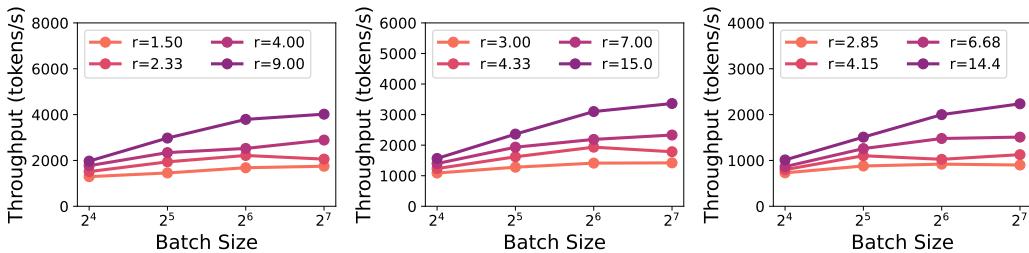


Figure 12: **MLP-to-Attention ratio on Inference Throughput (Qwen3):** (left) Qwen3-0.6B model variants; (center) Qwen3-1.7B model variants; (right) Qwen3-4B model variants. Across varying batch sizes and model scales, a larger MLP-to-Attention ratio increases inference throughput under a fixed parameter budget. The legend indicates the MLP-to-Attention ratio of the models, where  $r = r_{\text{mlp/attn}}$ . All evaluations are performed using the vLLM framework Kwon et al. (2023) on a single NVIDIA Ampere 40GB A100 GPU with 4096 input and 1024 output tokens.

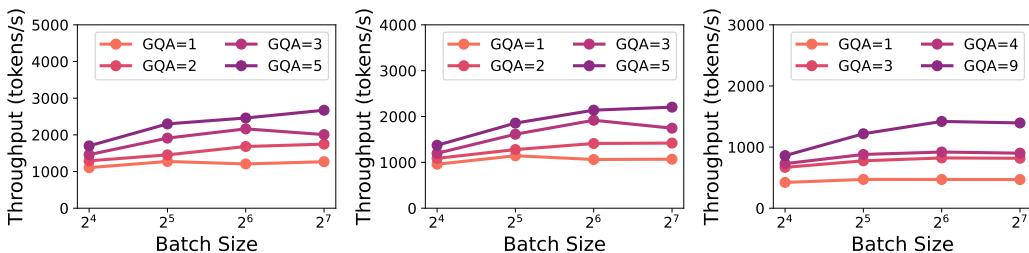


Figure 13: **GQA on Inference Throughput (Qwen3):** (left) Qwen3-0.6B model variants; (center) Qwen3-1.7B model variants; (right) Qwen3-4B model variants. This figure shows the impact of GQA on inference throughput. With the total parameter count fixed, hidden size is set to 1024 (0.6B), 2048 (1.7B), and 2560 (4B), and the MLP-to-Attention ratio is 1.5, 3.0, and 2.85, respectively. Across varying batch sizes, models with larger GQA achieve higher throughput. All evaluations are performed using the vLLM framework Kwon et al. (2023) on a single NVIDIA Ampere 40GB A100 GPU with 4096 input and 1024 output tokens.

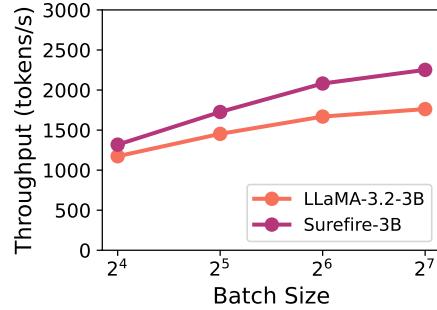
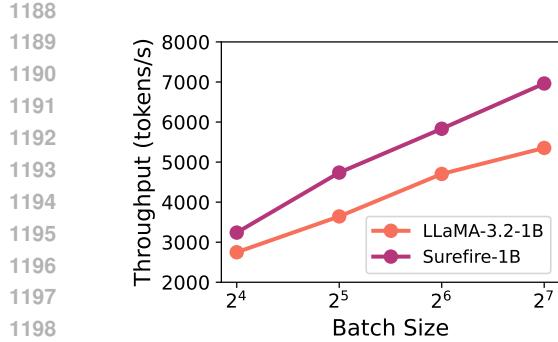


Figure 14: **Results for 1B and 3B models:** (left) Inference throughput comparison between LLaMA-3.2-1B and Surefire-1B, showing that Surefire-1B consistently achieves higher efficiency across batch sizes. (right) Inference throughput comparison between LLaMA-3.2-3B and Surefire-3B, demonstrating that Surefire-3B consistently delivers higher efficiency across all batch sizes. The results are collected using the SGLang framework Zheng et al. (2023) on a single A100 GPU with 4096 input and 1024 output tokens.

## F ADDITIONAL INFERENCE EVALUATION RESULTS OVER H200 GPUs

In this section, we present additional inference efficiency results on NVIDIA H200 GPUs. We derive architectural variants by altering the configurations of Qwen3-0.6B, 1.7B, and 4B to investigate the impact of model architectural factors on inference efficiency. The results are shown in Figure 15-17. We also compare the inference throughput of Surefire-1B, Surefire-3B, with LLaMA-3.2-1B and LLaMA-3.2-3B over H200 GPUs in Figure 18.

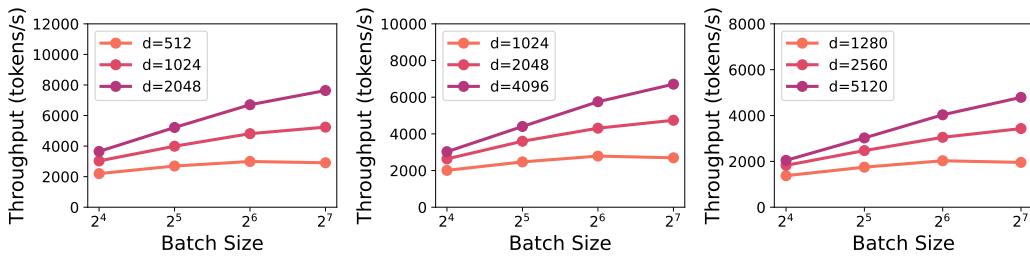


Figure 15: **Hidden size on Inference Throughput (Qwen3):** (left) Qwen3-0.6B model variants; (center) Qwen3-1.7B model variants; (right) Qwen3-4B model variants. Across varying batch sizes and model scales, larger hidden sizes yield higher inference throughput under a fixed parameter budget. The legend indicates the hidden size of the models, where  $d = d_{\text{model}}$ . All evaluations are performed using the vLLM framework Kwon et al. (2023) on a single NVIDIA H200 GPU with 4096 input and 1024 output tokens.

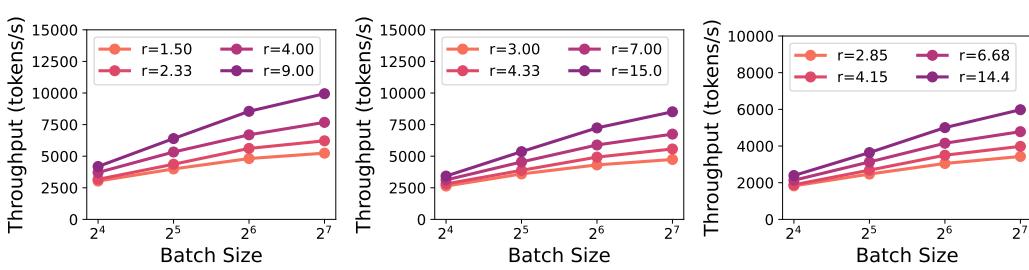


Figure 16: **MLP-to-Attention ratio on Inference Throughput (Qwen3):** (left) Qwen3-0.6B model variants; (center) Qwen3-1.7B model variants; (right) Qwen3-4B model variants. Across varying batch sizes and model scales, a larger MLP-to-Attention ratio increases inference throughput under a fixed parameter budget. The legend indicates the MLP-to-Attention ratio of the models, where  $r = r_{\text{mlp}/\text{attn}}$ . All evaluations are performed using the vLLM framework Kwon et al. (2023) on a single NVIDIA H200 GPU with 4096 input and 1024 output tokens.

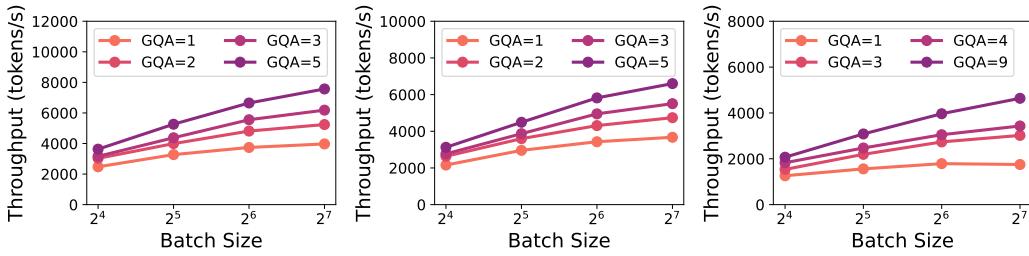


Figure 17: **GQA on Inference Throughput (Qwen3):** (left) Qwen3-0.6B model variants; (center) Qwen3-1.7B model variants; (right) Qwen3-4B model variants. This figure shows the impact of GQA on inference throughput. With the total parameter count fixed, hidden size is set to 1024 (0.6B), 2048 (1.7B), and 2560 (4B), and the MLP-to-Attention ratio is 1.5, 3.0, and 2.85, respectively. Across varying batch sizes, models with larger GQA achieve higher throughput. All evaluations are performed using the vLLM framework Kwon et al. (2023) on a single NVIDIA H200 GPU with 4096 input and 1024 output tokens.

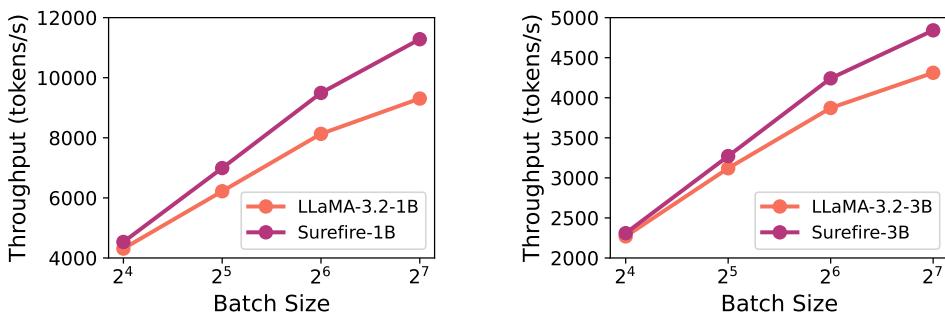
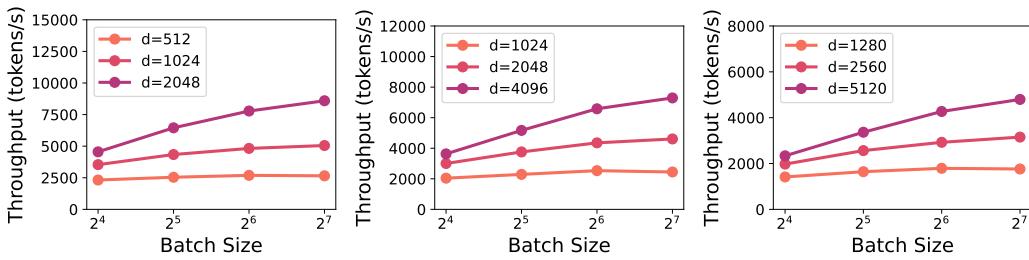


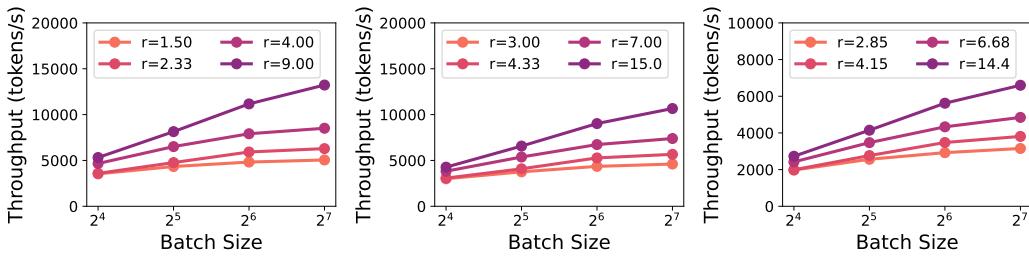
Figure 18: **Results for 1B and 3B models:** (left) Inference throughput comparison between LLaMA-3.2-1B and Surefire-1B, showing that Surefire-1B consistently achieves higher efficiency across batch sizes. (right) Inference throughput comparison between LLaMA-3.2-3B and Surefire-3B, demonstrating that Surefire-3B consistently delivers higher efficiency across all batch sizes. The results are collected using the SGLang framework Zheng et al. (2023) on a single NVIDIA H200 GPU with 4096 input and 1024 output tokens.

Furthermore, we use the SGLang framework Zheng et al. (2023) to measure the inference throughput of large language models. We construct architectural variants by modifying the configurations of

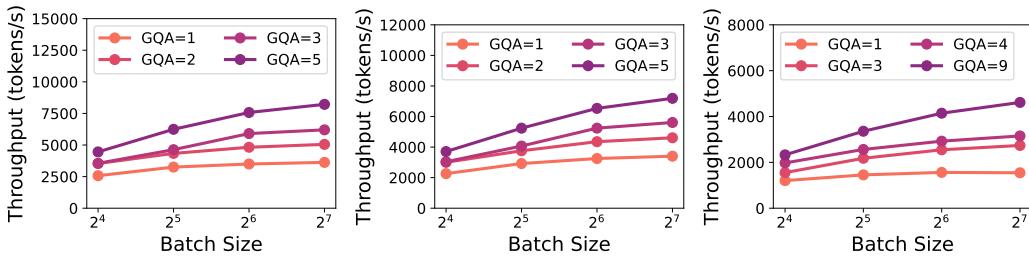
1296 Qwen3-0.6B, 1.7B, and 4B to study how different architectural factors influence inference efficiency.  
 1297 The results are presented in Figure 19-21. The inference throughput of Surefire-1B and Surefire-3B  
 1298 compared with LLaMA-3.2-1B and LLaMA-3.2-3B on H200 GPUs is shown in Figure 22.  
 1299



1300  
 1301  
 1302  
 1303  
 1304  
 1305  
 1306  
 1307  
 1308  
 1309  
**Figure 19: Hidden size on Inference Throughput (Qwen3):** (left) Qwen3-0.6B model variants;  
 1310 (center) Qwen3-1.7B model variants; (right) Qwen3-4B model variants. Across varying batch sizes  
 1311 and model scales, larger hidden sizes yield higher inference throughput under a fixed parameter  
 1312 budget. The legend indicates the hidden size of the models, where  $d = d_{\text{model}}$ . All evaluations are  
 1313 performed using the SGLang framework Zheng et al. (2023) on a single NVIDIA H200 GPU with  
 1314 4096 input and 1024 output tokens.  
 1315



1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322  
 1323  
 1324  
 1325  
 1326  
**Figure 20: MLP-to-Attention ratio on Inference Throughput (Qwen3):** (left) Qwen3-0.6B model  
 1327 variants; (center) Qwen3-1.7B model variants; (right) Qwen3-4B model variants. Across varying  
 1328 batch sizes and model scales, a larger MLP-to-Attention ratio increases inference throughput under  
 1329 a fixed parameter budget. The legend indicates the MLP-to-Attention ratio of the models, where  
 1330  $r = r_{\text{mlp/attn}}$ . All evaluations are performed using the SGLang framework Zheng et al. (2023) on a  
 1331 single NVIDIA H200 GPU with 4096 input and 1024 output tokens.  
 1332



1333  
 1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
**Figure 21: GQA on Inference Throughput (Qwen3):** (left) Qwen3-0.6B model variants; (center)  
 1344 Qwen3-1.7B model variants; (right) Qwen3-4B model variants. This figure shows the impact  
 1345 of GQA on inference throughput. With the total parameter count fixed, hidden size is set to 1024  
 1346 (0.6B), 2048 (1.7B), and 2560 (4B), and the MLP-to-Attention ratio is 1.5, 3.0, and 2.85, respec-  
 1347 tively. Across varying batch sizes, models with larger GQA achieve higher throughput. All eval-  
 1348 uations are performed using the SGLang framework Zheng et al. (2023) on a single NVIDIA H200  
 1349 GPU with 4096 input and 1024 output tokens.

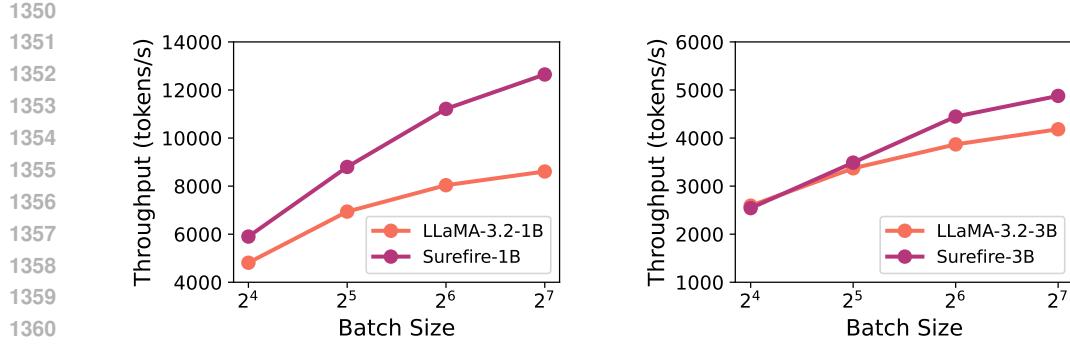


Figure 22: **Results for 1B and 3B models:** (left) Inference throughput comparison between LLaMA-3.2-1B and Surefire-1B, showing that Surefire-1B consistently achieves higher efficiency across batch sizes. (right) Inference throughput comparison between LLaMA-3.2-3B and Surefire-3B, demonstrating that Surefire-3B consistently delivers higher efficiency across all batch sizes. The results are collected using the SGLang framework Zheng et al. (2023) on a single NVIDIA H200 GPU with 4096 input and 1024 output tokens.

## G DETAILED THROUGHPUT STATISTICS

In this section, we present the detailed inference throughput results for the 1B and 3B models in Table 6.

Table 6: **Summary of Results for 1B and 3B Models:** We summarize the inference throughput (tokens/s) of LLaMA-3.2-1B, Surefire-1B, LLaMA-3.2-3B, and Surefire-3B across vLLM and SGLang on A100 and H200 GPUs using 4096 input tokens and 1024 output tokens.

| Hardware | Framework | Model        | Batch Size |         |          |          |
|----------|-----------|--------------|------------|---------|----------|----------|
|          |           |              | 16         | 32      | 64       | 128      |
| A100     | vLLM      | LLaMA-3.2-1B | 1931.87    | 2602.72 | 3409.85  | 3825.91  |
|          |           | Surefire-1B  | 2116.49    | 3290.23 | 4028.69  | 4800.05  |
|          |           | LLaMA-3.2-3B | 904.83     | 1121.39 | 1136.61  | 1222.03  |
|          |           | Surefire-3B  | 1005.44    | 1356.07 | 1613.32  | 1476.22  |
| A100     | SGLang    | LLaMA-3.2-1B | 2748.84    | 3643.27 | 4703.92  | 5353.29  |
|          |           | Surefire-1B  | 3239.55    | 4737.63 | 5832.01  | 6962.24  |
|          |           | LLaMA-3.2-3B | 1173.51    | 1452.97 | 1668.67  | 1762.18  |
|          |           | Surefire-3B  | 1318.23    | 1726.20 | 2081.44  | 2251.74  |
| H200     | vLLM      | LLaMA-3.2-1B | 4311.97    | 6221.14 | 8131.65  | 9306.36  |
|          |           | Surefire-1B  | 4532.85    | 6992.71 | 9493.46  | 11282.56 |
|          |           | LLaMA-3.2-3B | 2269.53    | 3119.94 | 3872.14  | 4311.43  |
|          |           | Surefire-3B  | 2309.48    | 3271.63 | 4242.33  | 4841.53  |
| H200     | SGLang    | LLaMA-3.2-1B | 4812.67    | 6939.88 | 8038.34  | 8608.57  |
|          |           | Surefire-1B  | 5900.52    | 8798.68 | 11214.40 | 12645.55 |
|          |           | LLaMA-3.2-3B | 2593.04    | 3370.42 | 3868.42  | 4183.09  |
|          |           | Surefire-3B  | 2542.21    | 3488.79 | 4446.66  | 4877.16  |

We further compare design choices across existing open-source models at the 1B and 3B scales in Table 7 and Table 8. For the LLaMA-3.2-1B, Panda-1B, and Surefire-1B models we pretrained, we report inference throughput (tokens/s), byte-level WikiText perplexity, and full architectural configurations in the accompanying tables. All throughput measurements are performed with vLLM on H200 GPUs using batch size 128. For the 1B scale, we include LLaMA-3.2-1B-HF and OLMo-2-1B-HF. Because OLMo supports only a 4k context window and cannot run our standard 4k/1k setup (4096 input tokens and 1024 output tokens), we additionally report results under a 2k/1k setup

(2048 input tokens and 1024 output tokens). For the 3B scale, we add LLaMA-3.2-3B-HF and Qwen2.5-3B-HF, all evaluated under the 4k/1k configuration.

**Table 7: Comparison against open-source models at the 1B scale:** We compare our pretrained LLaMA-3.2-1B, Panda-1B, and Surefire-1B models with LLaMA-3.2-1B-HF and OLMo-2-1B-HF in terms of inference throughput (on H200 GPUs using vLLM) and byte-level WikiText perplexity.

| Model                      | LLaMA-3.2-1B | Panda-1B | Surefire-1B | LLaMA-3.2-1B-HF | OLMo-2-1B-HF |
|----------------------------|--------------|----------|-------------|-----------------|--------------|
| Wikitext PPL               | 1.7151       | 1.7016   | 1.7142      | 1.5807          | 1.5798       |
| Tput (4k/1k)               | 9306         | 6218     | 11283       | 9306            | /            |
| Tput (2k/1k)               | 11948        | 8961     | 13890       | 11948           | 7486         |
| Model Architectural Config |              |          |             |                 |              |
| $n_{\text{layers}}$        | 16           | 16       | 16          | 16              | 16           |
| $d_{\text{model}}$         | 2048         | 2560     | 2560        | 2048            | 2048         |
| $r_{\text{mlp/attn}}$      | 4.8          | 1.067    | 3.6         | 4.8             | 3            |
| GQA                        | 4            | 4        | 9           | 4               | 1            |
| $N_{\text{non-embed}}$     | 973M         | 975M     | 965M        | 973M            | 1.074B       |

**Table 8: Comparison against open-source models at the 3B scale:** We compare our pretrained LLaMA-3.2-3B, Panda-3B, and Surefire-3B models with LLaMA-3.2-3B-HF and Qwen2.5-3B-HF in terms of inference throughput (on H200 GPUs using vLLM) and byte-level WikiText perplexity.

| Model                      | LLaMA-3.2-3B | Panda-3B | Surefire-3B | LLaMA-3.2-3B-HF | Qwen2.5-3B-HF |
|----------------------------|--------------|----------|-------------|-----------------|---------------|
| Wikitext PPL               | 1.6489       | 1.6454   | 1.6462      | 1.5164          | 1.6185        |
| Tput (4k/1k)               | 4311         | 3335     | 4842        | 4311            | 6470          |
| Model Architectural Config |              |          |             |                 |               |
| $n_{\text{layers}}$        | 28           | 28       | 28          | 28              | 36            |
| $d_{\text{model}}$         | 3072         | 4096     | 4096        | 3072            | 2048          |
| $r_{\text{mlp/attn}}$      | 3            | 1        | 1           | 3               | 7.167         |
| GQA                        | 3            | 3        | 7           | 3               | 8             |
| $N_{\text{non-embed}}$     | 2.82B        | 2.82B    | 2.82B       | 2.82B           | 2.77B         |

Our observations are as follows:

- OLMo-2-1B-HF is relatively close to our predicted optimal design, with an MLP-to-attention ratio of 3 (near our predicted 3.6), but remains inference-inefficient due to its hidden dimension and GQA choices.
- At the 3B scale, LLaMA-3.2-3B-HF achieves good accuracy but is not inference-efficient, while Qwen2.5-3B-HF is inference-efficient but less accurate.

These comparisons further underscore the necessity and relevance of our inference-efficient, high-accuracy model designs.

## 1458 H ADDITIONAL RESULTS: LOSS VS. GQA

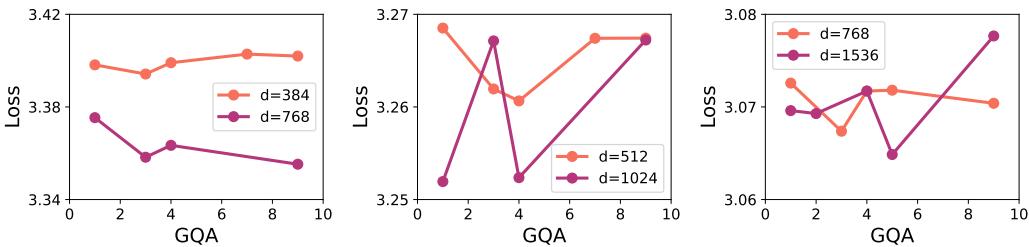
1460 We analyze the relationship between training loss and GQA while fixing the number of parameters,  
 1461 hidden size, and MLP-to-Attention ratio. In order to keep the number of attention parameters  
 1462 fixed, we vary GQA by holding the total number of heads fixed (i.e., total heads = query-heads +  
 1463 KV-heads) and re-allocating this fixed budget between query and key-value heads, producing asymmetric  
 1464 changes in their effective dimensionalities.

1465 As shown in Figure 23, unlike hidden size and MLP-to-Attention ratio, the relationship between loss  
 1466 and GQA is highly fluctuating. Varying GQA does not adjust model capacity in the coordinated way  
 1467 that changing  $d_{\text{model}}$  or  $r_{\text{mlp/attn}}$  does, where the dimensions of query, key, and value scale together  
 1468 predictably. Specifically, note the following facts when the total number of heads is fixed

- 1470 • Increasing the number of query-heads expands the query projection dimensionality but  
 1471 simultaneously reduces the number of KV-heads, increasing KV sharing and thus reducing  
 1472 KV expressivity.
- 1473 • Conversely, decreasing query-heads increases KV-head capacity (fewer replicas) but re-  
 1474 duces the projection dimensionality of both query and KV.

1475 These opposing effects create a tradeoff, making the relationship between GQA and training loss  
 1476 non-smooth and often highly fluctuating.

1477 Prior work shows only that query and KV projections can have non-interchangeable roles (e.g.,  
 1478 head-importance heterogeneity Voita et al. (2019)), but provides no monotonic or predictive theory  
 1479 for how reallocating capacity across query versus KV should affect loss. Consistent with this, recent  
 1480 open LLMs choose different GQA settings even within a single family: Qwen3 uses GQA = 2 for  
 1481 0.6B/1.7B, GQA = 4 for 4B/8B, GQA = 5 for 14B, GQA = 8 for 32B and for the 30B-A3B MoE;  
 1482 LLaMA-3/3.1/3.2 likewise use GQA = 4, 8, and 3 across closely related sizes. This variation across  
 1483 models of similar architecture shows that GQA is treated as a discrete, model-specific hyperparameter,  
 1484 supporting our decision to tune it via local search rather than integrate it into the continuous  
 1485 scaling law.



1487  
 1488 Figure 23: **Loss vs. GQA:** (left) 80M model variants; (center) 145M model variants; (right) 297M  
 1489 model variants. Across different model sizes, the relationship between training loss and GQA varies  
 1490 substantially when hidden size and the mlp-to-attention ratio are fixed. The legend denotes the  
 1491 hidden size of each trained model.

1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511

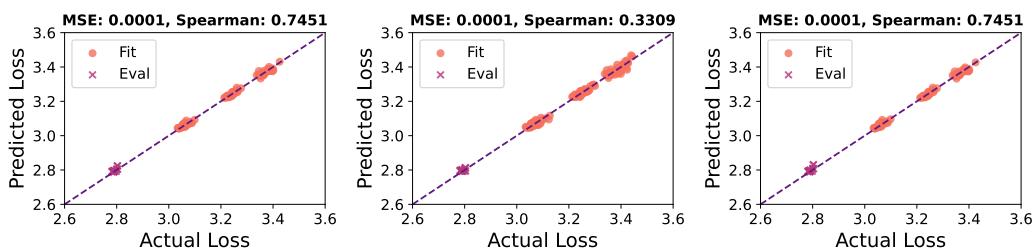
1512 **I MORE ABLATION STUDY**  
 1513

1514 In this section, We first evaluate the impact of outlier data on the fitting of the scaling laws in  
 1515 Figure 24 (left) and Figure 24 (center). Then, we evaluate the fitting performance of multiplicative  
 1516 calibrations and additive calibrations in Figure 24 (left) and Figure 24 (right).  
 1517

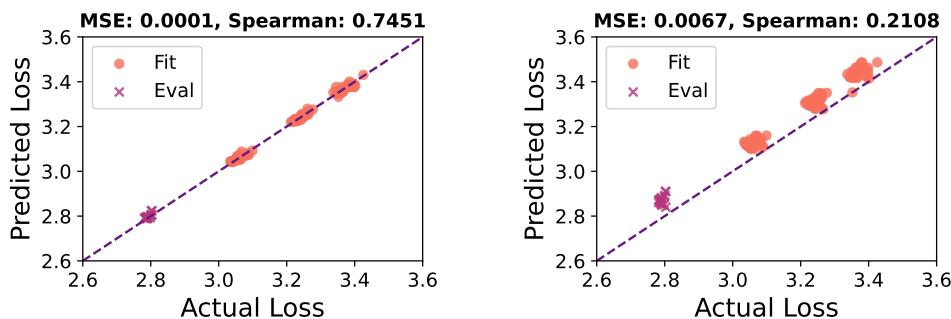
1518 Finally, we evaluate the performance of Joint and non-separable calibrations shown below in Figure 25:  
 1519

$$(a_0 + a_1 \log(\frac{dr}{\sqrt{N}}) + a_2 / (\frac{dr}{\sqrt{N}})) \cdot L_{\text{opt}}$$

1522 where  $d = d_{\text{model}}$ ,  $r = r_{\text{mlp/attn}}$ , and  $N = N_{\text{non-embed}}$ . In Figure 25, we observe that the performance  
 1523 of joint and non-separable calibrations is significantly worse than that of multiplicative calibration,  
 1524 consistent with our discussion in §3.3.  
 1525



1527 **Figure 24: Ablation Study:** (left) use multiplicative calibrations without outliers; (center) use  
 1528 multiplicative calibrations with outliers; (right) use additive calibrations without outliers. The outlier  
 1529 refers to models trained with an mlp-to-attention ratio below 0.5 or above 5. We observe that outlier  
 1530 data points harm the scaling law fit. Moreover, while multiplicative and additive calibrations differ  
 1531 in formulation, their MSE and Spearman values remain nearly identical. Dots denote the data points  
 1532 used for fitting, while crosses indicate the test data points.  
 1533



1535 **Figure 25: Joint and non-separable calibrations:** (left) use multiplicative calibrations; (right) use  
 1536 joint and non-separable calibrations. We observe that joint and non-separable calibrations yield  
 1537 higher MSE and lower Spearman scores than multiplicative calibrations, indicating inferior performance.  
 1538 Dots denote the data points used for fitting, while crosses indicate the test data points.  
 1539

1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

1566 **J INFERENCE FLOPs ANALYSIS**

1568 Building on the inference FLOPs analysis from prior work Kaplan et al. (2020), we begin with the  
 1569 following definition:

- 1570
- 1571 •  $d_{\text{model}}$ : hidden size
  - 1572 •  $f_{\text{size}}$ : intermediate (feed-forward) size
  - 1573 •  $n_{\text{layers}}$ : number of layers
  - 1574 •  $A$ : number of query heads
  - 1575 •  $K$ : number of key/value heads
  - 1576 •  $d_h$ : per-head hidden dimension (query and value)
  - 1577 •  $T$ : per-head hidden dim the KV length prior to token generation

1580 Based on the above definition, we have  $d_q = Ad_h$  and  $d_{kv} = Kd_h$ . We focus exclusively on  
 1581 non-embedding FLOPs, resulting in:

1582 Attention: QKV and Project

1583

$$n_{\text{layers}} \left( \underbrace{2d_{\text{model}}d_q}_Q + \underbrace{2d_{\text{model}}d_{kv}}_K + \underbrace{2d_{\text{model}}d_{kv}}_V + \underbrace{2d_{\text{model}}d_q}_O \right)$$

1587 Attention: Mask

1588

$$n_{\text{layers}} (2Td_q)$$

1590 Feedforward:

1591

$$n_{\text{layers}} (3 \cdot 2d_{\text{model}}f_{\text{size}})$$

1594 Total Inference non-embedding FLOPs:

1595

$$\text{Total-FLOPs} = n_{\text{layers}} \left( \underbrace{2d_{\text{model}}d_q}_Q + \underbrace{2d_{\text{model}}d_{kv}}_K + \underbrace{2d_{\text{model}}d_{kv}}_V + \underbrace{2d_{\text{model}}d_q}_O + \underbrace{2Td_q}_{qK^\top} + \underbrace{3 \cdot 2d_{\text{model}}f_{\text{size}}}_{\text{up, gate, down}} \right)$$

1599 Since  $P_{\text{non-emb}} \approx n_{\text{layers}}(2d_{\text{model}}d_q + 2d_{\text{model}}d_{kv} + 3d_{\text{model}}f_{\text{size}})$ . Therefore, Total-FLOPs =  
 1600  $2P_{\text{non-emb}} + 2n_{\text{layers}}Td_q$

1601 we adopt the following three approaches to accelerate inference:

- 1602
- 1603 • Increasing the MLP-to-Attention ratio reduces the term  $2Td_q$ , thereby lowering the total  
 1604 FLOPs.
  - 1605 • Increasing the hidden size reduces the term  $2Td_q$ , thereby lowering the total FLOPs.

## K MORE LARGE-SCALE TRAINING RESULTS

In this section, we first show the detailed result over downstream tasks of large-scale models in Table 9 and Table 10.

**Table 9: Detailed Results on Downstream Tasks for 1B Models:** In this table, we show detailed results of 1B models over 9 downstream tasks.

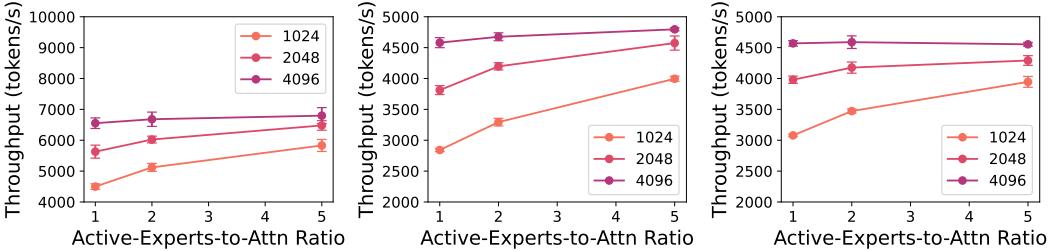
| Downstream Tasks | LLaMA-3.2-1B | Panda-1B | Surefire-1B |
|------------------|--------------|----------|-------------|
| Arc-Easy         | 58.8         | 60.9     | 59.7        |
| Arc-Challenge    | 29.8         | 28.9     | 30.2        |
| LAMBADA          | 52.8         | 55.1     | 52.0        |
| HellaSwag        | 56.9         | 58.4     | 56.6        |
| OpenBookQA       | 32.0         | 33.2     | 32.0        |
| PIQA             | 73.6         | 75.2     | 73.0        |
| SciQ             | 84.8         | 87.2     | 84.9        |
| WinoGrande       | 57.1         | 58.6     | 57.5        |
| COQA             | 48.7         | 55.3     | 52.7        |
| Avg.             | 54.9         | 57.0     | 55.4        |

**Table 10: Detailed Results on Downstream Tasks for 3B Models:** In this table, we show detailed results of 3B models over 9 downstream tasks.

| Downstream Tasks | LLaMA-3.2-3B | Panda-3B | Surefire-3B | Panda-3B° |
|------------------|--------------|----------|-------------|-----------|
| Arc-Easy         | 66.4         | 65.5     | 67.6        | 66.8      |
| Arc-Challenge    | 33.3         | 35.2     | 33.9        | 33.3      |
| LAMBADA          | 60.6         | 61.8     | 61.4        | 61.5      |
| HellaSwag        | 66.7         | 66.9     | 67.0        | 67.8      |
| OpenBookQA       | 38.4         | 38.6     | 38.6        | 38.0      |
| PIQA             | 76.8         | 76.9     | 77.4        | 76.8      |
| SciQ             | 89.4         | 91.2     | 92.1        | 90.5      |
| WinoGrande       | 62.5         | 63.2     | 60.5        | 62.7      |
| COQA             | 63.3         | 63.4     | 65.4        | 64.9      |
| Avg.             | 61.9         | 62.5     | 62.6        | 62.5      |

## 1674 L MOE INFERENCE

1675  
 1676 In this section, we examine how the Mixture-of-Experts (MoE) architecture affects inference effi-  
 1677 ciency. Figure 26 indicates that larger hidden sizes and higher Active-Experts-to-Attention ratios  
 1678 improve the inference throughput of MoE models, consistent with observations in dense models.  
 1679



1680  
 1681 Figure 26: **Active-Experts-to-Attn on Inference Throughput:** (left) 3B-A1.1B model variants;  
 1682 (center) 5.3B-A1.7B model variants; (right) 8.3B-A1.5B model variants. We study the effect of  
 1683 the Active-Experts-to-Attention ratio on inference throughput by fixing the total number of active  
 1684 parameters, setting GQA to 4, and using a batch size of 2048 to reduce MoE inference variance in  
 1685 this figure. All evaluations are performed using the vLLM framework Kwon et al. (2023) on a single  
 1686 NVIDIA Ampere 40GB A100 GPU with 1024 input and 256 output tokens.  
 1687  
 1688

1689  
 1690  
 1691  
 1692  
 1693  
 1694  
 1695  
 1696  
 1697  
 1698  
 1699  
 1700  
 1701  
 1702  
 1703  
 1704  
 1705  
 1706  
 1707  
 1708  
 1709  
 1710  
 1711  
 1712  
 1713  
 1714  
 1715  
 1716  
 1717  
 1718  
 1719  
 1720  
 1721  
 1722  
 1723  
 1724  
 1725  
 1726  
 1727